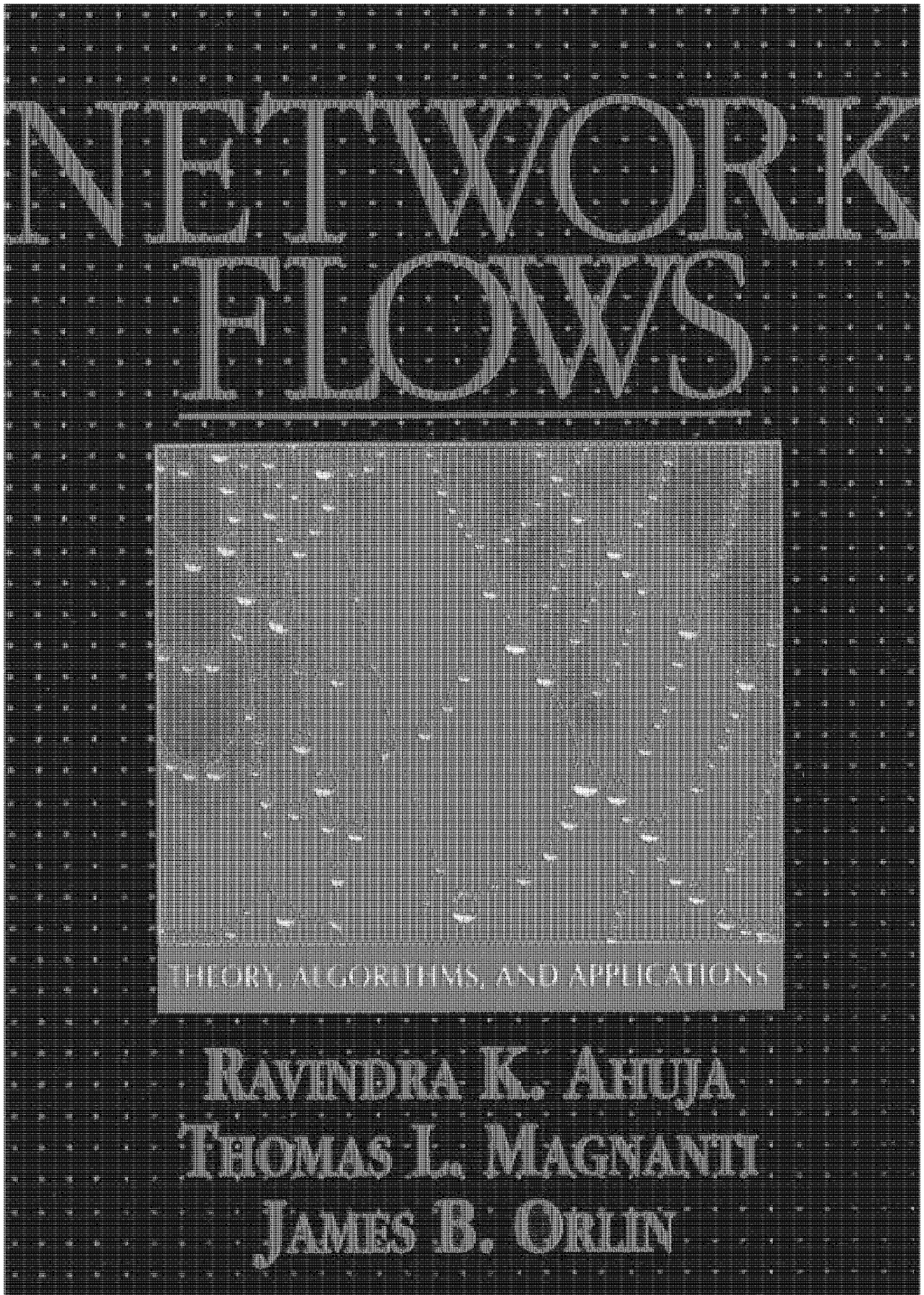


EXHIBIT 4



NETWORK FLOWS

Theory, Algorithms, and Applications

RAVINDRA K. AHUJA

Department of Industrial & Management Engineering
Indian Institute of Technology, Kanpur

THOMAS L. MAGNANTI

Sloan School of Management
Massachusetts Institute of Technology, Cambridge

JAMES B. ORLIN

Sloan School of Management
Massachusetts Institute of Technology, Cambridge



PRENTICE HALL, Upper Saddle River, New Jersey 07458

Library of Congress Cataloging-in-Publication Data

Ahuja, Ravindra K. (date)
Network flows : theory, algorithms, and applications / Ravindra K.
Ahuja, Thomas L. Magnanti, James B. Orlin.
p. cm.
Includes bibliographical references and index.
ISBN 0-13-617549-X
1. Network analysis (Planning) 2. Mathematical optimization.
I. Magnanti, Thomas L. II. Orlin, James B., (date). III. Title.
T57.85.A37 1993
658.4'032—dc20

92-26702
CIP

Acquisitions editor: Pete Janzow
Production editor: Merrill Peterson
Cover designer: Design Source
Prepress buyer: Linda Behrens
Manufacturing buyer: David Dickey
Editorial assistant: Phyllis Morgan

The author and publisher of this book have used their best efforts in preparing this book. These efforts include the development, research, and testing of the theories and programs to determine their effectiveness. The author and publisher make no warranty of any kind, expressed or implied, with regard to these programs or the documentation contained in this book. The author and publisher shall not be liable in any event for incidental or consequential damages in connection with, or arising out of, the furnishing, performance, or use of these programs.



© 1993 by Prentice-Hall, Inc.
Upper Saddle River, New Jersey 07458

All rights reserved. No part of this book may be reproduced, in any form or by any means, without permission in writing from the publisher.

Printed in the United States of America

16 17 18 19

ISBN 0-13-617549-X

PRENTICE-HALL INTERNATIONAL (UK) LIMITED, *London*
PRENTICE-HALL OF AUSTRALIA PTY. LIMITED, *Sydney*
PRENTICE-HALL CANADA INC., *Toronto*
PRENTICE-HALL HISPANOAMERICANA, S.A., *Mexico*
PRENTICE-HALL OF INDIA PRIVATE LIMITED, *New Delhi*
PRENTICE-HALL OF JAPAN, INC., *Tokyo*
EDITORA PRENTICE-HALL DO BRASIL, LTDA., *Rio de Janeiro*

CONTENTS

PREFACE, xi

1 INTRODUCTION, 1

- 1.1 Introduction, 1
- 1.2 Network Flow Problems, 4
- 1.3 Applications, 9
- 1.4 Summary, 18
- Reference Notes, 19
- Exercises, 20

2 PATHS, TREES, AND CYCLES, 23

- 2.1 Introduction, 23
- 2.2 Notation and Definitions, 24
- 2.3 Network Representations, 31
- 2.4 Network Transformations, 38
- 2.5 Summary, 46
- Reference Notes, 47
- Exercises, 47

3 ALGORITHM DESIGN AND ANALYSIS, 53

- 3.1 Introduction, 53
- 3.2 Complexity Analysis, 56
- 3.3 Developing Polynomial-Time Algorithms, 66
- 3.4 Search Algorithms, 73
- 3.5 Flow Decomposition Algorithms, 79
- 3.6 Summary, 84
- Reference Notes, 85
- Exercises, 86

4 SHORTEST PATHS: LABEL-SETTING ALGORITHMS, 93

- 4.1 Introduction, 93
- 4.2 Applications, 97
- 4.3 Tree of Shortest Paths, 106
- 4.4 Shortest Path Problems in Acyclic Networks, 107
- 4.5 Dijkstra's Algorithm, 108
- 4.6 Dial's Implementation, 113
- 4.7 Heap Implementations, 115
- 4.8 Radix Heap Implementation, 116

4.9	Summary, 121
	Reference Notes, 122
	Exercises, 124
5	SHORTEST PATHS: LABEL-CORRECTING ALGORITHMS, 133
5.1	Introduction, 133
5.2	Optimality Conditions, 135
5.3	Generic Label-Correcting Algorithms, 136
5.4	Special Implementations of the Modified Label-Correcting Algorithm, 141
5.5	Detecting Negative Cycles, 143
5.6	All-Pairs Shortest Path Problem, 144
5.7	Minimum Cost-to-Time Ratio Cycle Problem, 150
5.8	Summary, 154
	Reference Notes, 156
	Exercises, 157
6	MAXIMUM FLOWS: BASIC IDEAS, 166
6.1	Introduction, 166
6.2	Applications, 169
6.3	Flows and Cuts, 177
6.4	Generic Augmenting Path Algorithm, 180
6.5	Labeling Algorithm and the Max-Flow Min-Cut Theorem, 184
6.6	Combinatorial Implications of the Max-Flow Min-Cut Theorem, 188
6.7	Flows with Lower Bounds, 191
6.8	Summary, 196
	Reference Notes, 197
	Exercises, 198
7	MAXIMUM FLOWS: POLYNOMIAL ALGORITHMS, 207
7.1	Introduction, 207
7.2	Distance Labels, 209
7.3	Capacity Scaling Algorithm, 210
7.4	Shortest Augmenting Path Algorithm, 213
7.5	Distance Labels and Layered Networks, 221
7.6	Generic Preflow-Push Algorithm, 223
7.7	FIFO Preflow-Push Algorithm, 231
7.8	Highest-Label Preflow-Push Algorithm, 233
7.9	Excess Scaling Algorithm, 237
7.10	Summary, 241
	Reference Notes, 241
	Exercises, 243
8	MAXIMUM FLOWS: ADDITIONAL TOPICS, 250
8.1	Introduction, 250
8.2	Flows in Unit Capacity Networks, 252
8.3	Flows in Bipartite Networks, 255
8.4	Flows in Planar Undirected Networks, 260
8.5	Dynamic Tree Implementations, 265

- 8.6 Network Connectivity, 273
- 8.7 All-Pairs Minimum Value Cut Problem, 277
- 8.8 Summary, 285
- Reference Notes, 287
- Exercises, 288

9 MINIMUM COST FLOWS: BASIC ALGORITHMS, 294

- 9.1 Introduction, 294
- 9.2 Applications, 298
- 9.3 Optimality Conditions, 306
- 9.4 Minimum Cost Flow Duality, 310
- 9.5 Relating Optimal Flows to Optimal Node Potentials, 315
- 9.6 Cycle-Canceling Algorithm and the Integrality Property, 317
- 9.7 Successive Shortest Path Algorithm, 320
- 9.8 Primal-Dual Algorithm, 324
- 9.9 Out-of-Kilter Algorithm, 326
- 9.10 Relaxation Algorithm, 332
- 9.11 Sensitivity Analysis, 337
- 9.12 Summary, 339
- Reference Notes, 341
- Exercises, 344

10 MINIMUM COST FLOWS: POLYNOMIAL ALGORITHMS, 357

- 10.1 Introduction, 357
- 10.2 Capacity Scaling Algorithm, 360
- 10.3 Cost Scaling Algorithm, 362
- 10.4 Double Scaling Algorithm, 373
- 10.5 Minimum Mean Cycle-Canceling Algorithm, 376
- 10.6 Repeated Capacity Scaling Algorithm, 382
- 10.7 Enhanced Capacity Scaling Algorithm, 387
- 10.8 Summary, 395
- Reference Notes, 396
- Exercises, 397

11 MINIMUM COST FLOWS: NETWORK SIMPLEX ALGORITHMS, 402

- 11.1 Introduction, 402
- 11.2 Cycle Free and Spanning Tree Solutions, 405
- 11.3 Maintaining a Spanning Tree Structure, 409
- 11.4 Computing Node Potentials and Flows, 411
- 11.5 Network Simplex Algorithm, 415
- 11.6 Strongly Feasible Spanning Trees, 421
- 11.7 Network Simplex Algorithm for the Shortest Path Problem, 425
- 11.8 Network Simplex Algorithm for the Maximum Flow Problem, 430
- 11.9 Related Network Simplex Algorithms, 433
- 11.10 Sensitivity Analysis, 439
- 11.11 Relationship to Simplex Method, 441
- 11.12 Unimodularity Property, 447
- 11.13 Summary, 450
- Reference Notes, 451
- Exercises, 453

12 ASSIGNMENTS AND MATCHINGS, 461

- 12.1 Introduction, 461
- 12.2 Applications, 463
- 12.3 Bipartite Cardinality Matching Problem, 469
- 12.4 Bipartite Weighted Matching Problem, 470
- 12.5 Stable Marriage Problem, 473
- 12.6 Nonbipartite Cardinality Matching Problem, 475
- 12.7 Matchings and Paths, 494
- 12.8 Summary, 498
- Reference Notes, 499
- Exercises, 501

13 MINIMUM SPANNING TREES, 510

- 13.1 Introduction, 510
- 13.2 Applications, 512
- 13.3 Optimality Conditions, 516
- 13.4 Kruskal's Algorithm, 520
- 13.5 Prim's Algorithm, 523
- 13.6 Sollin's Algorithm, 526
- 13.7 Minimum Spanning Trees and Matroids, 528
- 13.8 Minimum Spanning Trees and Linear Programming, 530
- 13.9 Summary, 533
- Reference Notes, 535
- Exercises, 536

14 CONVEX COST FLOWS, 543

- 14.1 Introduction, 543
- 14.2 Applications, 546
- 14.3 Transformation to a Minimum Cost Flow Problem, 551
- 14.4 Pseudopolynomial-Time Algorithms, 554
- 14.5 Polynomial-Time Algorithm, 556
- 14.6 Summary, 560
- Reference Notes, 561
- Exercises, 562

15 GENERALIZED FLOWS, 566

- 15.1 Introduction, 566
- 15.2 Applications, 568
- 15.3 Augmented Forest Structures, 572
- 15.4 Determining Potentials and Flows for an Augmented Forest Structure, 577
- 15.5 Good Augmented Forests and Linear Programming Bases, 582
- 15.6 Generalized Network Simplex Algorithm, 583
- 15.7 Summary, 591
- Reference Notes, 591
- Exercises, 593

16 LAGRANGIAN RELAXATION AND NETWORK OPTIMIZATION, 598

- 16.1 Introduction, 598
- 16.2 Problem Relaxations and Branch and Bound, 602
- 16.3 Lagrangian Relaxation Technique, 605
- 16.4 Lagrangian Relaxation and Linear Programming, 615
- 16.5 Applications of Lagrangian Relaxation, 620
- 16.6 Summary, 635
- Reference Notes, 637
- Exercises, 638

17 MULTICOMMODITY FLOWS, 649

- 17.1 Introduction, 649
- 17.2 Applications, 653
- 17.3 Optimality Conditions, 657
- 17.4 Lagrangian Relaxation, 660
- 17.5 Column Generation Approach, 665
- 17.6 Dantzig-Wolfe Decomposition, 671
- 17.7 Resource-Directive Decomposition, 674
- 17.8 Basis Partitioning, 678
- 17.9 Summary, 682
- Reference Notes, 684
- Exercises, 686

18 COMPUTATIONAL TESTING OF ALGORITHMS, 695

- 18.1 Introduction, 695
- 18.2 Representative Operation Counts, 698
- 18.3 Application to Network Simplex Algorithm, 702
- 18.4 Summary, 713
- Reference Notes, 713
- Exercises, 715

19 ADDITIONAL APPLICATIONS, 717

- 19.1 Introduction, 717
- 19.2 Maximum Weight Closure of a Graph, 719
- 19.3 Data Scaling, 725
- 19.4 Science Applications, 728
- 19.5 Project Management, 732
- 19.6 Dynamic Flows, 737
- 19.7 Arc Routing Problems, 740
- 19.8 Facility Layout and Location, 744
- 19.9 Production and Inventory Planning, 748
- 19.10 Summary, 755
- Reference Notes, 759
- Exercises, 760

APPENDIX A DATA STRUCTURES, 765

- A.1 Introduction, 765
- A.2 Elementary Data Structures, 766
- A.3 d -Heaps, 773
- A.4 Fibonacci Heaps, 779
- Reference Notes, 787

APPENDIX B \mathcal{NP} -COMPLETENESS, 788

- B.1 Introduction, 788
- B.2 Problem Reductions and Transformations, 790
- B.3 Problem Classes \mathcal{P} , \mathcal{NP} , \mathcal{NP} -Complete, and \mathcal{NP} -Hard, 792
- B.4 Proving \mathcal{NP} -Completeness Results, 796
- B.5 Concluding Remarks, 800
- Reference Notes, 801

APPENDIX C LINEAR PROGRAMMING, 802

- C.1 Introduction, 802
- C.2 Graphical Solution Procedure, 804
- C.3 Basic Feasible Solutions, 805
- C.4 Simplex Method, 810
- C.5 Bounded Variable Simplex Method, 814
- C.6 Linear Programming Duality, 816
- Reference Notes, 820

REFERENCES, 821

INDEX, 840

6

MAXIMUM FLOWS: BASIC IDEAS

*You get the maxx for the minimum at T. J. Maxx.**
—Advertisement for a clothing store.

Chapter Outline

-
- 6.1 Introduction
 - 6.2 Applications
 - 6.3 Flows and Cuts
 - 6.4 Generic Augmenting Path Algorithm
 - 6.5 Labeling Algorithm and the Max-Flow Min-Cut Theorem
 - 6.6 Combinatorial Implications of the Max-Flow Min-Cut Theorem
 - 6.7 Flows with Lower Bounds
 - 6.8 Summary
-

6.1 INTRODUCTION

The maximum flow problem and the shortest path problem are complementary. They are similar because they are both pervasive in practice and because they both arise as subproblems in algorithms for the minimum cost flow problem. The two problems differ, however, because they capture different aspects of the minimum cost flow problem: Shortest path problems model arc costs but not arc capacities; maximum flow problems model capacities but not costs. Taken together, the shortest path problem and the maximum flow problem combine all the basic ingredients of network flows. As such, they have become the nuclei of network optimization. Our study of the shortest path problem in the preceding two chapters has introduced us to some of the basic building blocks of network optimization, such as distance labels, optimality conditions, and some core strategies for designing iterative solution methods and for improving the performance of these methods. Our discussion of maximum flows, which we begin in this chapter, builds on these ideas and introduces several other key ideas that reoccur often in the study of network flows.

The *maximum flow problem* is very easy to state: In a capacitated network, we wish to send as much flow as possible between two special nodes, a source node s and a sink node t , without exceeding the capacity of any arc. In this and the following two chapters, we discuss a number of algorithms for solving the maximum flow problem. These algorithms are of two types:

* © The TJX Operating Companies, Inc. 1985. Get the Maxx for the Minimum® and T. J. Maxx are registered trademarks of the TJX Operating Companies, Inc.

1. Augmenting path algorithms that maintain mass balance constraints at every node of the network other than the source and sink nodes. These algorithms incrementally augment flow along paths from the source node to the sink node.
2. Preflow-push algorithms that flood the network so that some nodes have excesses (or buildup of flow). These algorithms incrementally relieve flow from nodes with excesses by sending flow from the node forward toward the sink node or backward toward the source node.

We discuss the simplest version of the first type of algorithm in this chapter and more elaborate algorithms of both types in Chapter 7. To help us to understand the importance of the maximum flow problem, we begin by describing several applications. This discussion shows how maximum flow problems arise in settings as diverse as manufacturing, communication systems, distribution planning, matrix rounding, and scheduling.

We begin our algorithmic discussion by considering a *generic augmenting path algorithm* for solving the maximum flow problem and describing an important special implementation of the generic approach, known as the *labeling algorithm*. The labeling algorithm is a pseudopolynomial-time algorithm. In Chapter 7 we develop improved versions of this generic approach with better theoretical behavior. The correctness of these algorithms rests on the renowned *max-flow min-cut theorem* of network flows (recall from Section 2.2 that a cut is a set of arcs whose deletion disconnects the network into two parts). This central theorem in the study of network flows (indeed, perhaps the most significant theorem in this problem domain) not only provides us with an instrument for analyzing algorithms, but also permits us to model a variety of applications in machine and vehicle scheduling, communication systems planning, and several other settings, as maximum flow problems, even though on the surface these problems do not appear to have a network flow structure. In Section 6.6 we describe several such applications.

The max-flow min-cut theorem establishes an important correspondence between flows and cuts in networks. Indeed, as we will see, by solving a maximum flow problem, we also solve a complementary *minimum cut problem*: From among all cuts in the network that separate the source and sink nodes, find the cut with the minimum capacity. The relationship between maximum flows and minimum cuts is important for several reasons. First, it embodies a fundamental duality result that arises in many problem settings in discrete mathematics and that underlies linear programming as well as mathematical optimization in general. In fact, the max-flow min-cut theorem, which shows the equivalence between the maximum flow and minimum cut problems, is a special case of the well-known strong duality theorem of linear programming. The fact that maximum flow problems and minimum cut problems are equivalent has practical implications as well. It means that the theory and algorithms that we develop for the maximum flow problem are also applicable to many practical problems that are naturally cast as minimum cut problems. Our discussion of combinatorial applications in the text and exercises of this chapter and our discussion of applications in Chapter 19 features several applications of this nature.

Notation and Assumptions

We consider a capacitated network $G = (N, A)$ with a *nonnegative* capacity u_{ij} associated with each arc $(i, j) \in A$. Let $U = \max\{u_{ij} : (i, j) \in A\}$. As before, the arc adjacency list $A(i) = \{(i, k) : (i, k) \in A\}$ contains all the arcs emanating from node i . To define the maximum flow problem, we distinguish two special nodes in the network G : a *source node* s and a *sink node* t . We wish to find the maximum flow from the source node s to the sink node t that satisfies the arc capacities and mass balance constraints at all nodes. We can state the problem formally as follows.

$$\text{Maximize } v \quad (6.1a)$$

subject to

$$\sum_{\{j : (i,j) \in A\}} x_{ij} - \sum_{\{j : (j,i) \in A\}} x_{ji} = \begin{cases} v & \text{for } i = s, \\ 0 & \text{for all } i \in N - \{s \text{ and } t\} \\ -v & \text{for } i = t \end{cases} \quad (6.1b)$$

$$0 \leq x_{ij} \leq u_{ij} \quad \text{for each } (i, j) \in A. \quad (6.1c)$$

We refer to a vector $x = \{x_{ij}\}$ satisfying (6.1b) and (6.1c) as a *flow* and the corresponding value of the scalar variable v as the *value* of the flow. We consider the maximum flow problem subject to the following assumptions.

Assumption 6.1. *The network is directed.*

As explained in Section 2.4, we can always fulfill this assumption by transforming any undirected network into a directed network.

Assumption 6.2. *All capacities are nonnegative integers.*

Although it is possible to relax the integrality assumption on arc capacities for some algorithms, this assumption is necessary for others. Algorithms whose complexity bounds involve U assume integrality of the data. In reality, the integrality assumption is not a restrictive assumption because all modern computers store capacities as rational numbers and we can always transform rational numbers to integer numbers by multiplying them by a suitably large number.

Assumption 6.3. *The network does not contain a directed path from node s to node t composed only of infinite capacity arcs.*

Whenever every arc on a directed path P from s to t has infinite capacity, we can send an infinite amount of flow along this path, and therefore the maximum flow value is unbounded. Notice that we can detect the presence of an infinite capacity path using the search algorithm described in Section 3.4.

Assumption 6.4. *Whenever an arc (i, j) belongs to A , arc (j, i) also belongs to A .*

This assumption is nonrestrictive because we allow arcs with zero capacity.

Assumption 6.5. *The network does not contain parallel arcs (i.e., two or more arcs with the same tail and head nodes).*

This assumption is essentially a notational convenience. In Exercise 6.24 we ask the reader to show that this assumption imposes no loss of generality.

Before considering the theory underlying the maximum flow problem and algorithms for solving it, and to provide some background and motivation for studying the problem, we first describe some applications.

6.2 APPLICATIONS

The maximum flow problem, and the minimum cut problem, arise in a wide variety of situations and in several forms. For example, sometimes the maximum flow problem occurs as a subproblem in the solution of more difficult network problems, such as the minimum cost flow problem or the generalized flow problem. As we will see in Section 6.6, the maximum flow problem also arises in a number of combinatorial applications that on the surface might not appear to be maximum flow problems at all. The problem also arises directly in problems as far reaching as machine scheduling, the assignment of computer modules to computer processors, the rounding of census data to retain the confidentiality of individual households, and tanker scheduling. In this section we describe a few such applications; in Chapter 19 we discuss several other applications.

Application 6.1 Feasible Flow Problem

The feasible flow problem requires that we identify a flow x in a network $G = (N, A)$ satisfying the following constraints:

$$\sum_{\{j:(i,j) \in A\}} x_{ij} - \sum_{\{j:(j,i) \in A\}} x_{ji} = b(i) \quad \text{for } i \in N, \quad (6.2a)$$

$$0 \leq x_{ij} \leq u_{ij} \quad \text{for all } (i, j) \in A. \quad (6.2b)$$

As before, we assume that $\sum_{i \in N} b(i) = 0$. The following distribution scenario illustrates how the feasible flow problem arises in practice. Suppose that merchandise is available at some seaports and is desired by other ports. We know the stock of merchandise available at the ports, the amount required at the other ports, and the maximum quantity of merchandise that can be shipped on a particular sea route. We wish to know whether we can satisfy all of the demands by using the available supplies.

We can solve the feasible flow problem by solving a maximum flow problem defined on an augmented network as follows. We introduce two new nodes, a source node s and a sink node t . For each node i with $b(i) > 0$, we add an arc (s, i) with capacity $b(i)$, and for each node i with $b(i) < 0$, we add an arc (i, t) with capacity $-b(i)$. We refer to the new network as the *transformed network*. Then we solve a maximum flow problem from node s to node t in the transformed network. If the maximum flow saturates all the source and sink arcs, problem (6.2) has a feasible solution; otherwise, it is infeasible. (In Section 6.7 we give necessary and sufficient conditions for a feasible flow problem to have a feasible solution.)

It is easy to verify why this algorithm works. If x is a flow satisfying (6.2a)

and (6.2b), the same flow with $x_{si} = b(i)$ for each source arc (s, i) and $x_{it} = -b(i)$ for each sink arc (i, t) is a maximum flow in the transformed network (since it saturates all the source and the sink arcs). Similarly, if x is a maximum flow in the transformed network that saturates all the source and the sink arcs, this flow in the original network satisfies (6.2a) and (6.2b). Therefore, the original network contains a feasible flow if and only if the transformed network contains a flow that saturates all the source and sink arcs. This observation shows how the maximum flow problem arises whenever we need to find a feasible solution in a network.

Application 6.2 Problem of Representatives

A town has r residents R_1, R_2, \dots, R_r ; q clubs C_1, C_2, \dots, C_q ; and p political parties P_1, P_2, \dots, P_p . Each resident is a member of at least one club and can belong to exactly one political party. Each club must nominate one of its members to represent it on the town's governing council so that the number of council members belonging to the political party P_k is at most u_k . Is it possible to find a council that satisfies this "balancing" property?

We illustrate this formulation with an example. We consider a problem with $r = 7$, $q = 4$, $p = 3$, and formulate it as a maximum flow problem in Figure 6.1. The nodes R_1, R_2, \dots, R_7 represent the residents, the nodes C_1, C_2, \dots, C_4 represent the clubs, and the nodes P_1, P_2, \dots, P_3 represent the political parties.

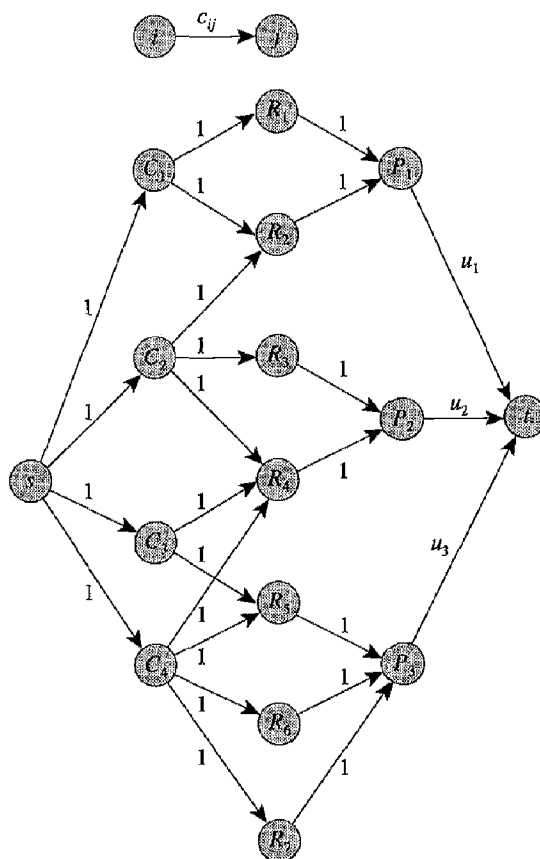


Figure 6.1 System of distinct representatives.

The network also contains a source node s and a sink node t . It contains an arc (s, C_i) for each node C_i denoting a club, an arc (C_i, R_j) whenever the resident R_j is a member of the club C_i , and an arc (R_j, P_k) if the resident R_j belongs to the political party P_k . Finally, we add an arc (P_k, t) for each $k = 1, \dots, 3$ of capacity u_k ; all other arcs have unit capacity.

We next find a maximum flow in this network. If the maximum flow value equals q , the town has a balanced council; otherwise, it does not. The proof of this assertion is easy to establish by showing that (1) any flow of value q in the network corresponds to a balanced council, and that (2) any balanced council implies a flow of value q in the network.

This type of model has applications in several resource assignment settings. For example, suppose that the residents are skilled craftsmen, the club C_i is the set of craftsmen with a particular skill, and the political party P_k corresponds to a particular seniority class. In this instance, a balanced town council corresponds to an assignment of craftsmen to a union governing board so that every skill class has representation on the board and no seniority class has a dominant representation.

Application 6.3 Matrix Rounding Problem

This application is concerned with consistent rounding of the elements, row sums, and column sums of a matrix. We are given a $p \times q$ matrix of *real* numbers $D = \{d_{ij}\}$, with row sums α_i and column sums β_j . We can round any real number a to the next smaller integer $\lfloor a \rfloor$ or to the next larger integer $\lceil a \rceil$, and the decision to round up or down is entirely up to us. The matrix rounding problem requires that we round the matrix elements, and the row and column sums of the matrix so that the sum of the rounded elements in each row equals the rounded row sum and the sum of the rounded elements in each column equals the rounded column sum. We refer to such a rounding as a *consistent rounding*.

We shall show how we can discover such a rounding scheme, if it exists, by solving a feasible flow problem for a network with nonnegative lower bounds on arc flows. (As shown in Section 6.7, we can solve this problem by solving two maximum flow problems with zero lower bounds on arc flows.) We illustrate our method using the matrix rounding problem shown in Figure 6.2. Figure 6.3 shows the maximum flow network for this problem. This network contains a node i corresponding to each row i and a node j' corresponding to each column j . Observe that this network

				Row sum
	3.1	6.8	7.3	17.2
	9.6	2.4	0.7	12.7
	3.6	1.2	6.5	11.3
Column sum	16.3	10.4	14.5	

Figure 6.2 Matrix rounding problem.

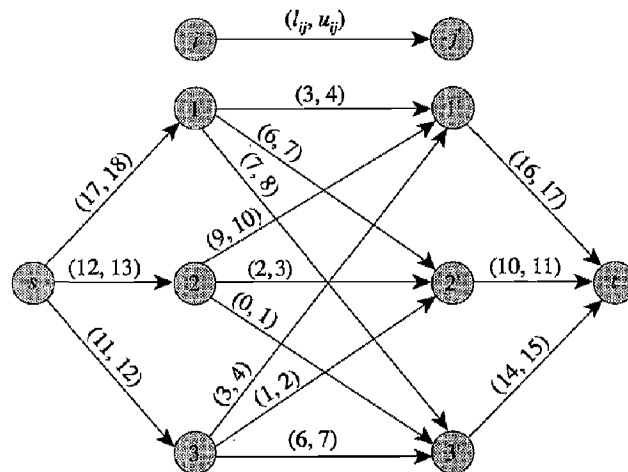


Figure 6.3 Network for the matrix rounding problem.

contains an arc (i, j') for each matrix element d_{ij} , an arc (s, i) for each row sum, and an arc (j', t) for each column sum. The lower and the upper bounds of each arc (i, j') are $\lfloor d_{ij} \rfloor$ and $\lceil d_{ij} \rceil$, respectively. It is easy to establish a one-to-one correspondence between the consistent roundings of the matrix and feasible flows in the corresponding network. Consequently, we can find a consistent rounding by solving a maximum flow problem on the corresponding network.

This matrix rounding problem arises in several application contexts. For example, the U.S. Census Bureau uses census information to construct millions of tables for a wide variety of purposes. By law, the bureau has an obligation to protect the source of its information and not disclose statistics that could be attributed to any particular person. We might disguise the information in a table as follows. We round off each entry in the table, including the row and column sums, either up or down to a multiple of a constant k (for some suitable value of k), so that the entries in the table continue to add to the (rounded) row and column sums, and the overall sum of the entries in the new table adds to a rounded version of the overall sums in the original table. This Census Bureau problem is the same as the matrix rounding problem discussed earlier except that we need to round each element to a multiple of $k \geq 1$ instead of rounding it to a multiple of 1. We solve this problem by defining the associated network as before, but now defining the lower and upper bounds for any arc with an associated real number α as the greatest multiple of k less than or equal to α and the smallest multiple of k greater than or equal to α .

Application 6.4 Scheduling on Uniform Parallel Machines

In this application we consider the problem of scheduling of a set J of jobs on M uniform parallel machines. Each job $j \in J$ has a processing requirement p_j (denoting the number of machine days required to complete the job), a release date r_j (representing the beginning of the day when job j becomes available for processing), and a due date $d_j \geq r_j + p_j$ (representing the beginning of the day by which the job must be completed). We assume that a machine can work on only one job at a time and that each job can be processed by at most one machine at a time. However, we

allow *preemptions* (i.e., we can interrupt a job and process it on different machines on different days). The scheduling problem is to determine a feasible schedule that completes all jobs before their due dates or to show that no such schedule exists.

Scheduling problems like this arise in batch processing systems involving batches with a large number of units. The feasible scheduling problem, described in the preceding paragraph, is a fundamental problem in this situation and can be used as a subroutine for more general scheduling problems, such as the maximum lateness problem, the (weighted) minimum completion time problem, and the (weighted) maximum utilization problem.

Let us formulate the feasible scheduling problem as a maximum flow problem. We illustrate the formulation using the scheduling problem described in Figure 6.4 with $M = 3$ machines. First, we rank all the release and due dates, r_j and d_j for all j , in ascending order and determine $P \leq 2|J| - 1$ mutually disjoint intervals of dates between consecutive milestones. Let $T_{k,l}$ denote the interval that starts at the beginning of date k and ends at the beginning of date $l + 1$. For our example, this order of release and due dates is 1, 3, 4, 5, 7, 9. We have five intervals, represented by $T_{1,2}$, $T_{3,3}$, $T_{4,4}$, $T_{5,6}$, and $T_{7,8}$. Notice that within each interval $T_{k,l}$, the set of available jobs (i.e., those released but not yet due) does not change: we can process all jobs j with $r_j \leq k$ and $d_j \geq l + 1$ in the interval.

Job (j)	1	2	3	4
Processing time (p_j)	1.5	1.25	2.1	3.6
Release time (r_j)	3	1	3	5
Due date (d_j)	5	4	7	9

Figure 6.4 Scheduling problem.

We formulate the scheduling problem as a maximum flow problem on a bipartite network G as follows. We introduce a source node s , a sink node t , a node corresponding to each job j , and a node corresponding to each interval $T_{k,l}$, as shown in Figure 6.5. We connect the source node to every job node j with an arc with capacity p_j , indicating that we need to assign p_j days of machine time to job j . We connect each interval node $T_{k,l}$ to the sink node t by an arc with capacity $(l - k + 1)M$, representing the total number of machine days available on the days from k to l . Finally, we connect a job node j to every interval node $T_{k,l}$ if $r_j \leq k$ and $d_j \geq l + 1$ by an arc with capacity $(l - k + 1)$ which represents the maximum number of machines days we can allot to job j on the days from k to l . We next solve a maximum flow problem on this network: The scheduling problem has a feasible schedule if and only if the maximum flow value equals $\sum_{j \in J} p_j$ [alternatively, the flow on every arc (s, j) is p_j]. The validity of this formulation is easy to establish by showing a one-to-one correspondence between feasible schedules and flows of value $\sum_{j \in J} p_j$ from the source to the sink.

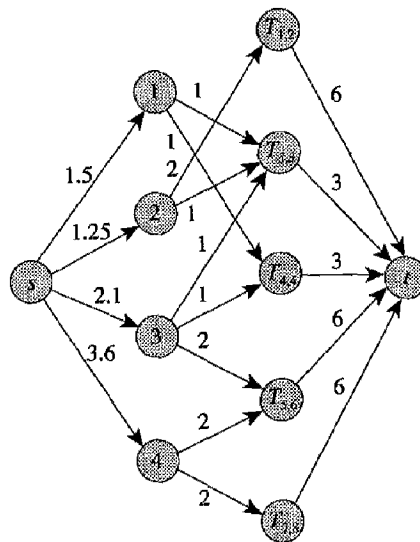


Figure 6.5 Network for scheduling uniform parallel machines.

Application 6.5 Distributed Computing on a Two-Processor Computer

This application concerns assigning different modules (subroutines) of a program to two processors in a way that minimizes the collective costs of interprocessor communication and computation. We consider a computer system with two processors; they need not be identical. We wish to execute a large program on this computer system. Each program contains several modules that interact with each other during the program's execution. The cost of executing each module on the two processes is known in advance and might vary from one processor to the other because of differences in the processors' memory, control, speed, and arithmetic capabilities. Let α_i and β_i denote the cost of computation of module i on processors 1 and 2, respectively. Assigning different modules to different processors incurs relatively high overhead costs due to interprocessor communication. Let c_{ij} denote the interprocessor communication cost if modules i and j are assigned to different processors; we do not incur this cost if we assign modules i and j to the same processor. The cost structure might suggest that we allocate two jobs to different processors—we need to balance this cost against the communication costs that we incur by allocating the jobs to different processors. Therefore, we wish to allocate modules of the program on the two processors so that we minimize the total cost of processing and interprocessor communication.

We formulate this problem as a minimum cut problem on an undirected network as follows. We define a source node s representing processor 1, a sink node t representing processor 2, and a node for every module of the program. For every node i , other than the source and sink nodes, we include an arc (s, i) of capacity β_i and an arc (i, t) of capacity α_i . Finally, if module i interacts with module j during program execution, we include the arc (i, j) with a capacity equal to c_{ij} . Figures 6.6 and 6.7 give an example of this construction. Figure 6.6 gives the data for this problem, and Figure 6.7 gives the corresponding network.

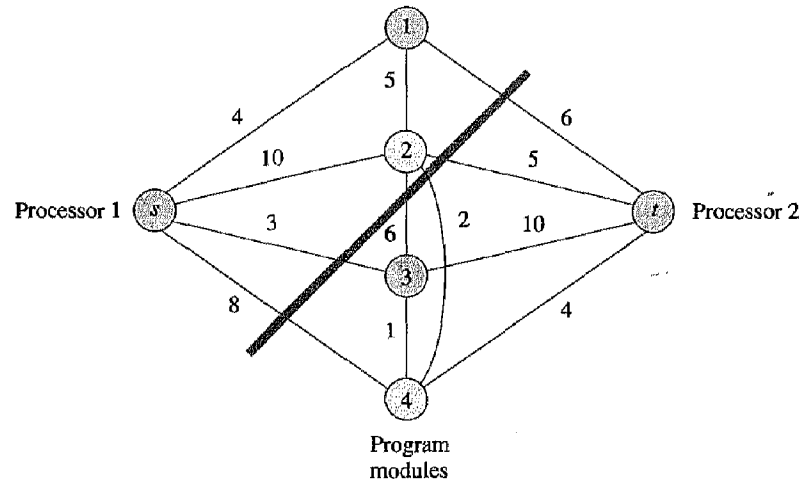
We now observe a one-to-one correspondence between s – t cuts in the network

i	1	2	3	4
α_i	6	5	10	4
β_i	4	10	3	8

(a)

	1	2	3	4
1	0	5	0	0
2	5	0	6	2
3	0	6	0	1
4	0	2	1	0

(b)

Figure 6.6 Data for the distributed computing model.**Figure 6.7** Network for the distributed computing model.

and assignments of modules to the two processors; moreover, the capacity of a cut equals the cost of the corresponding assignment. To establish this result, let A_1 and A_2 be an assignment of modules to processors 1 and 2, respectively. The cost of this assignment is $\sum_{i \in A_1} \alpha_i + \sum_{i \in A_2} \beta_i + \sum_{(i,j) \in A_1 \times A_2} c_{ij}$. The s - t cut corresponding to this assignment is $(\{s\} \cup A_1, \{t\} \cup A_2)$. The approach we used to construct the network implies that this cut contains an arc (i, t) for every $i \in A_1$ of capacity α_i , an arc (s, i) for every $i \in A_2$ of capacity β_i , and all arcs (i, j) with $i \in A_1$ and $j \in A_2$ with capacity c_{ij} . The cost of the assignment A_1 and A_2 equals the capacity of the cut $(\{s\} \cup A_1, \{t\} \cup A_2)$. (We suggest that readers verify this conclusion using

the example given in Figure 6.7 with $A_1 = \{1, 2\}$ and $A_2 = \{3, 4\}$.) Consequently, the minimum s - t cut in the network gives the minimum cost assignment of the modules to the two processors.

Application 6.6 Tanker Scheduling Problem

A steamship company has contracted to deliver perishable goods between several different origin–destination pairs. Since the cargo is perishable, the customers have specified precise dates (i.e., delivery dates) when the shipments must reach their destinations. (The cargoes may not arrive early or late.) The steamship company wants to determine the minimum number of ships needed to meet the delivery dates of the shiploads.

To illustrate a modeling approach for this problem, we consider an example with four shipments; each shipment is a full shipload with the characteristics shown in Figure 6.8(a). For example, as specified by the first row in this figure, the company must deliver one shipload available at port A and destined for port C on day 3. Figure 6.8(b) and (c) show the transit times for the shipments (including allowances for loading and unloading the ships) and the return times (without a cargo) between the ports.

Ship-ment	Origin	Desti-nation	Delivery date
1	Port A	Port C	3
2	Port A	Port C	8
3	Port B	Port D	3
4	Port B	Port C	6

(a)

	<i>C</i>	<i>D</i>
<i>A</i>	3	2
<i>B</i>	2	3

(b)

	<i>A</i>	<i>B</i>
<i>C</i>	2	1
<i>D</i>	1	2

(c)

Figure 6.8 Data for the tanker scheduling problem: (a) shipment characteristics; (b) shipment transit times; (c) return times.

We solve this problem by constructing a network shown in Figure 6.9(a). This network contains a node for each shipment and an arc from node i to node j if it is possible to deliver shipment j after completing shipment i ; that is, the start time of shipment j is no earlier than the delivery time of shipment i plus the travel time from the destination of shipment i to the origin of shipment j . A directed path in this network corresponds to a feasible sequence of shipment pickups and deliveries. The tanker scheduling problem requires that we identify the minimum number of directed paths that will contain each node in the network on exactly one path.

We can transform this problem to the framework of the maximum flow problem as follows. We split each node i into two nodes i' and i'' and add the arc (i', i'') . We set the lower bound on each arc (i', i'') , called the *shipment arc*, equal to 1 so that at least one unit of flow passes through this arc. We also add a source node s and connect it to the origin of each shipment (to represent putting a ship into service),

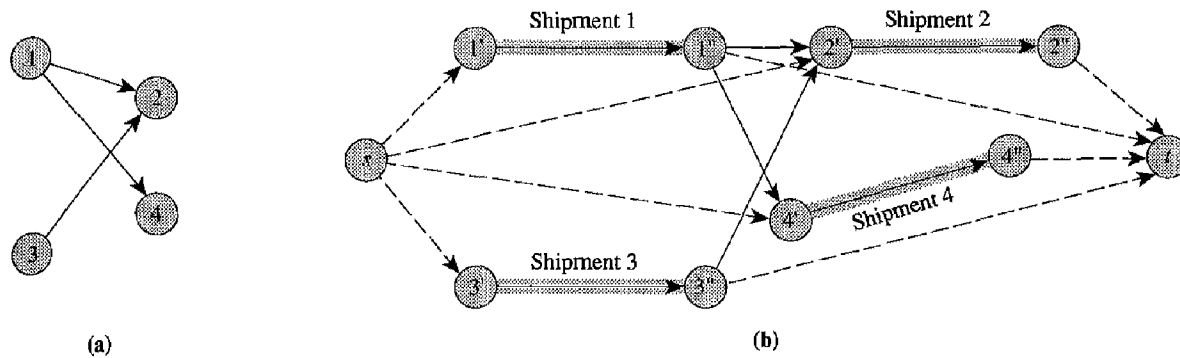


Figure 6.9 Network formulation of the tanker scheduling problem: (a) network of feasible sequences of two consecutive shipments; (b) maximum flow model.

and we add a sink node t and connect each destination node to it (to represent taking a ship out of service). We set the capacity of each arc in the network to value 1. Figure 6.9(b) shows the resulting network for our example. In this network, each directed path from the source s to the sink t corresponds to a feasible schedule for a single ship. As a result, a feasible flow of value v in this network decomposes into schedules of v ships and our problem reduces to identifying a feasible flow of minimum value. We note that the zero flow is not feasible because shipment arcs have unit lower bounds. We can solve this problem, which is known as the *minimum value problem*, using any maximum flow algorithm (see Exercise 6.18).

6.3 FLOWS AND CUTS

In this section we discuss some elementary properties of flows and cuts. We use these properties to prove the max-flow min-cut theorem to establish the correctness of the generic augmenting path algorithm. We first review some of our previous notation and introduce a few new ideas.

Residual network. The concept of *residual network* plays a central role in the development of all the maximum flow algorithms we consider. Earlier in Section 2.4 we defined residual networks and discussed several of its properties. Given a flow x , the residual capacity r_{ij} of any arc $(i, j) \in A$ is the maximum additional flow that can be sent from node i to node j using the arcs (i, j) and (j, i) . [Recall our assumption from Section 6.1 that whenever the network contains arc (i, j) , it also contains arc (j, i) .] The residual capacity r_{ij} has two components: (1) $u_{ij} - x_{ij}$, the unused capacity of arc (i, j) , and (2) the current flow x_{ji} on arc (j, i) , which we can cancel to increase the flow from node i to node j . Consequently, $r_{ij} = u_{ij} - x_{ij} + x_{ji}$. We refer to the network $G(x)$ consisting of the arcs with positive residual capacities as the *residual network* (with respect to the flow x). Figure 6.10 gives an example of a residual network.

s - t cut. We now review notation about cuts. Recall from Section 2.2 that a cut is a partition of the node set N into two subsets S and $\bar{S} = N - S$; we represent this cut using the notation $[S, \bar{S}]$. Alternatively, we can define a cut as the set of

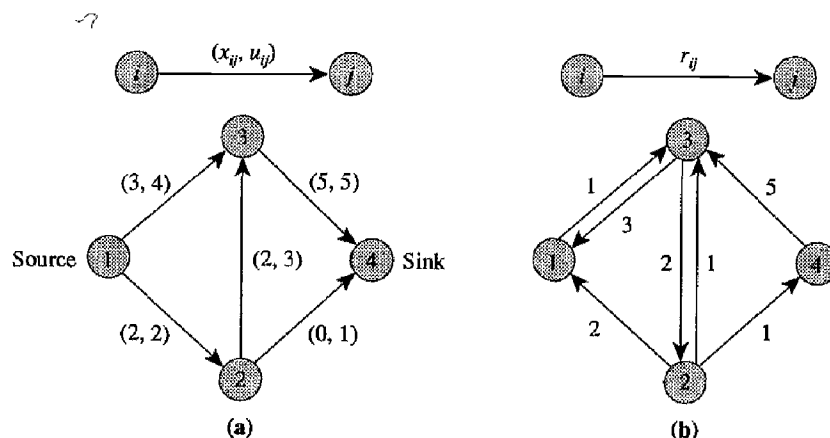


Figure 6.10 Illustrating a residual network: (a) original network G with a flow x ; (b) residual network $G(x)$.

arcs whose endpoints belong to the different subsets S and \bar{S} . We refer to a cut as an s - t cut if $s \in S$ and $t \in \bar{S}$. We also refer to an arc (i, j) with $i \in S$ and $j \in \bar{S}$ as a *forward arc* of the cut, and an arc (i, j) with $i \in \bar{S}$ and $j \in S$ as a *backward arc* of the cut $[S, \bar{S}]$. Let (S, \bar{S}) denote the set of forward arcs in the cut, and let (\bar{S}, S) denote the set of backward arcs. For example, in Figure 6.11, the dashed arcs constitute an s - t cut. For this cut, $(S, \bar{S}) = \{(1, 2), (3, 4), (5, 6)\}$, and $(\bar{S}, S) = \{(2, 3), (4, 5)\}$.

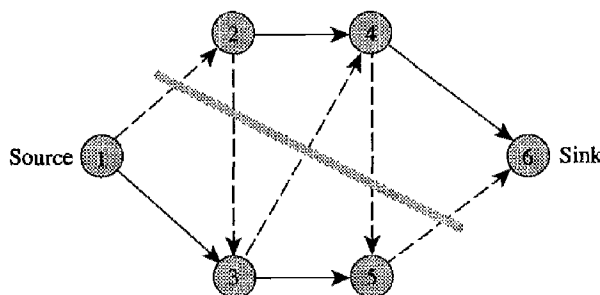


Figure 6.11 Example of an s - t cut.

Capacity of an s - t cut. We define the capacity $u[S, \bar{S}]$ of an s - t cut $[S, \bar{S}]$ as the sum of the capacities of the forward arcs in the cut. That is,

$$u[S, \bar{S}] = \sum_{(i,j) \in (S, \bar{S})} u_{ij}.$$

Clearly, the capacity of a cut is an upper bound on the maximum amount of flow we can send from the nodes in S to the nodes in \bar{S} while honoring arc flow bounds.

Minimum cut. We refer to an s - t cut whose capacity is minimum among all s - t cuts as a *minimum cut*.

Residual capacity of an s - t cut. We define the residual capacity $r[S, \bar{S}]$ of an s - t cut $[S, \bar{S}]$ as the sum of the residual capacities of forward arcs in the cut. That is,

$$r[S, \bar{S}] = \sum_{(i,j) \in (S, \bar{S})} r_{ij}.$$

Flow across an s - t cut. Let x be a flow in the network. Adding the mass balance constraint (6.1b) for the nodes in S , we see that

$$v = \sum_{i \in S} \left[\sum_{\{j: (i,j) \in A\}} x_{ij} - \sum_{\{j: (j,i) \in A\}} x_{ji} \right].$$

We can simplify this expression by noting that whenever both the nodes p and q belong to S and $(p, q) \in A$, the variable x_{pq} in the first term within the brackets (for node $i = p$) cancels the variable $-x_{pq}$ in the second term within the brackets (for node $j = q$). Moreover, if both the nodes p and q belong to \bar{S} , then x_{pq} does not appear in the expression. This observation implies that

$$v = \sum_{(i,j) \in (S, \bar{S})} x_{ij} - \sum_{(i,j) \in (\bar{S}, S)} x_{ij}. \quad (6.3)$$

The first expression on the right-hand side of (6.3) denotes the amount of flow from the nodes in S to nodes in \bar{S} , and the second expression denotes the amount of flow returning from the nodes in \bar{S} to the nodes in S . Therefore, the right-hand side denotes the total (net) flow across the cut, and (6.3) implies that the flow across any s - t cut $[S, \bar{S}]$ equals v . Substituting $x_{ij} \leq u_{ij}$ in the first expression of (6.3) and $x_{ij} \geq 0$ in the second expression shows that

$$v \leq \sum_{(i,j) \in (S, \bar{S})} u_{ij} = u[S, \bar{S}]. \quad (6.4)$$

This expression indicates that the value of any flow is less than or equal to the capacity of any s - t cut in the network. This result is also quite intuitive. Any flow from node s to node t must pass through every s - t cut in the network (because any cut divides the network into two disjoint components), and therefore the value of the flow can never exceed the capacity of the cut. Let us formally record this result.

Property 6.1. *The value of any flow is less than or equal to the capacity of any cut in the network.*

This property implies that if we discover a flow x whose value equals the capacity of some cut $[S, \bar{S}]$, then x is a maximum flow and the cut $[S, \bar{S}]$ is a minimum cut. The max-flow min-cut theorem, proved in the next section, states that some flow always has a flow value equal to the capacity of some cut.

We next restate Property 6.1 in terms of the residual capacities. Suppose that x is a flow of value v . Moreover, suppose that x' is a flow of value $v + \Delta v$ for some $\Delta v \geq 0$. The inequality (6.4) implies that

$$v + \Delta v \leq \sum_{(i,j) \in (S, \bar{S})} u_{ij}. \quad (6.5)$$

Subtracting (6.3) from (6.5) shows that

$$\Delta v \leq \sum_{(i,j) \in (S, \bar{S})} (u_{ij} - x_{ij}) + \sum_{(i,j) \in (\bar{S}, S)} x_{ij}. \quad (6.6)$$

We now use Assumption 6.4 to note that we can rewrite $\sum_{(i,j) \in (\bar{S}, S)} x_{ij}$ as $\sum_{(i,j) \in (S, \bar{S})} x_{ji}$. Consequently,

$$\Delta v \leq \sum_{(i,j) \in (S, \bar{S})} (u_{ij} - x_{ij} + x_{ji}) = \sum_{(S, \bar{S})} r_{ij}.$$

The following property is now immediate.

Property 6.2. *For any flow x of value v in a network, the additional flow that can be sent from the source node s to the sink node t is less than or equal to the residual capacity of any s – t cut.*

6.4 GENERIC AUGMENTING PATH ALGORITHM

In this section, we describe one of the simplest and most intuitive algorithms for solving the maximum flow problem. This algorithm is known as the *augmenting path algorithm*.

We refer to a directed path from the source to the sink in the residual network as an *augmenting path*. We define the residual capacity of an augmenting path as the minimum residual capacity of any arc in the path. For example, the residual network in Figure 6.10(b), contains exactly one augmenting path 1–3–2–4, and the residual capacity of this path is $\delta = \min\{r_{13}, r_{32}, r_{24}\} = \min\{1, 2, 1\} = 1$. Observe that, by definition, the capacity δ of an augmenting path is always positive. Consequently, whenever the network contains an augmenting path, we can send additional flow from the source to the sink. The generic augmenting path algorithm is essentially based on this simple observation. The algorithm proceeds by identifying augmenting paths and augmenting flows on these paths until the network contains no such path. Figure 6.12 describes the generic augmenting path algorithm.

```

algorithm augmenting path;
begin
   $x := 0$ ;
  while  $G(x)$  contains a directed path from node  $s$  to node  $t$  do
    begin
      identify an augmenting path  $P$  from node  $s$  to node  $t$ ;
       $\delta := \min\{r_{ij} : (i, j) \in P\}$ ;
      augment  $\delta$  units of flow along  $P$  and update  $G(x)$ ;
    end;
  end;

```

Figure 6.12 Generic augmenting path algorithm.

We use the maximum flow problem given in Figure 6.13(a) to illustrate the algorithm. Suppose that the algorithm selects the path 1–3–4 for augmentation. The residual capacity of this path is $\delta = \min\{r_{13}, r_{34}\} = \min\{4, 5\} = 4$. This augmentation reduces the residual capacity of arc (1, 3) to zero (thus we delete it from the residual network) and increases the residual capacity of arc (3, 1) to 4 (so we add this arc to the residual network). The augmentation also decreases the residual capacity of arc (3, 4) from 5 to 1 and increases the residual capacity of arc (4, 3) from 0 to 4. Figure 6.13(b) shows the residual network at this stage. In the second iteration, suppose that the algorithm selects the path 1–2–3–4. The residual capacity of this

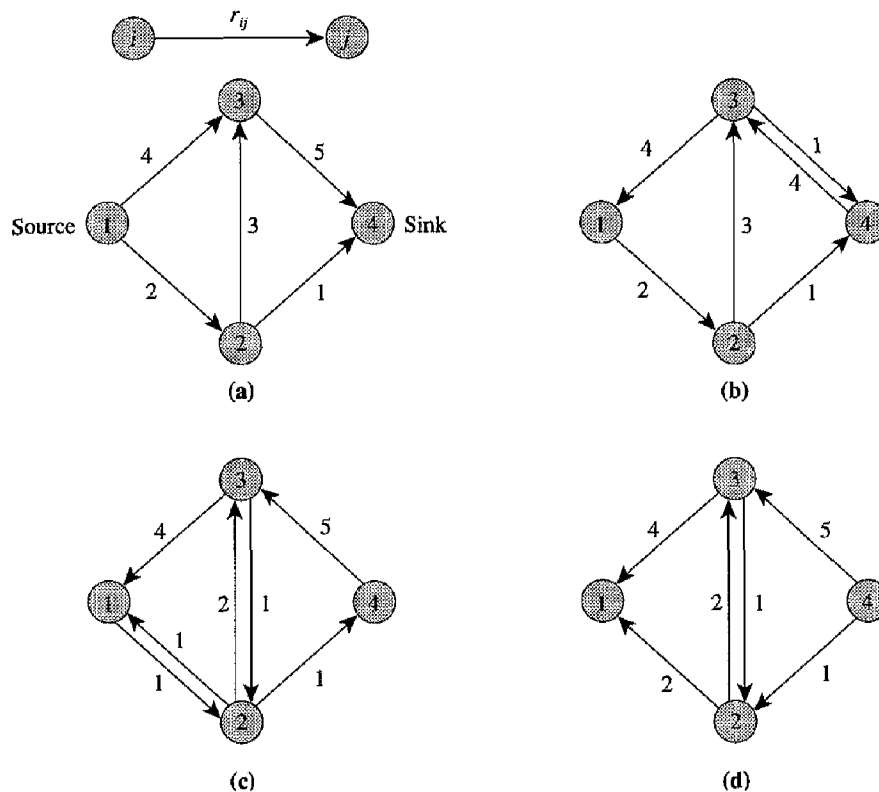


Figure 6.13 Illustrating the generic augmenting path algorithm: (a) residual network for the zero flow; (b) network after augmenting four units along the path 1-3-4; (c) network after augmenting one unit along the path 1-2-3-4; (d) network after augmenting one unit along the path 1-2-4.

path is $\delta = \min\{2, 3, 1\} = 1$. Augmenting 1 unit of flow along this path yields the residual network shown in Figure 6.13(c). In the third iteration, the algorithm augments 1 unit of flow along the path 1-2-4. Figure 6.13(d) shows the corresponding residual network. Now the residual network contains no augmenting path, so the algorithm terminates.

Relationship between the Original and Residual Networks

In implementing any version of the generic augmenting path algorithm, we have the option of working directly on the original network with the flows x_{ij} , or maintaining the residual network $G(x)$ and keeping track of the residual capacities r_{ij} and, when the algorithm terminates, recovering the actual flow variables x_{ij} . To see how we can use either alternative, it is helpful to understand the relationship between arc flows in the original network and residual capacities in the residual network.

First, let us consider the concept of an augmenting path in the original network. An augmenting path in the original network G is a path P (not necessarily directed) from the source to the sink with $x_{ij} < u_{ij}$ on every forward arc (i, j) and $x_{ij} > 0$ on every backward arc (i, j) . It is easy to show that the original network G contains

an augmenting path with respect to a flow x if and only if the residual network $G(x)$ contains a directed path from the source to the sink.

Now suppose that we update the residual capacities at some point in the algorithm. What is the effect on the arc flows x_{ij} ? The definition of the residual capacity (i.e., $r_{ij} = u_{ij} - x_{ij} + x_{ji}$) implies that an additional flow of δ units on arc (i, j) in the residual network corresponds to (1) an increase in x_{ij} by δ units in the original network, or (2) a decrease in x_{ji} by δ units in the original network, or (3) a convex combination of (1) and (2). We use the example given in Figure 6.14(a) and the corresponding residual network in Figure 6.14(b) to illustrate these possibilities. Augmenting 1 unit of flow on the path 1–2–4–3–5–6 in the network produces the residual network in Figure 6.14(c) with the corresponding arc flows shown in Figure 6.14(d). Comparing the solution in Figure 6.14(d) with that in Figure 6.14(a), we find that the flow augmentation increases the flow on arcs (1, 2), (2, 4), (3, 5), (5, 6) and decreases the flow on arc (3, 4).

Finally, suppose that we are given values for the residual capacities. How should we determine the flows x_{ij} ? Observe that since $r_{ij} = u_{ij} - x_{ij} + x_{ji}$, many combinations of x_{ij} and x_{ji} correspond to the same value of r_{ij} . We can determine

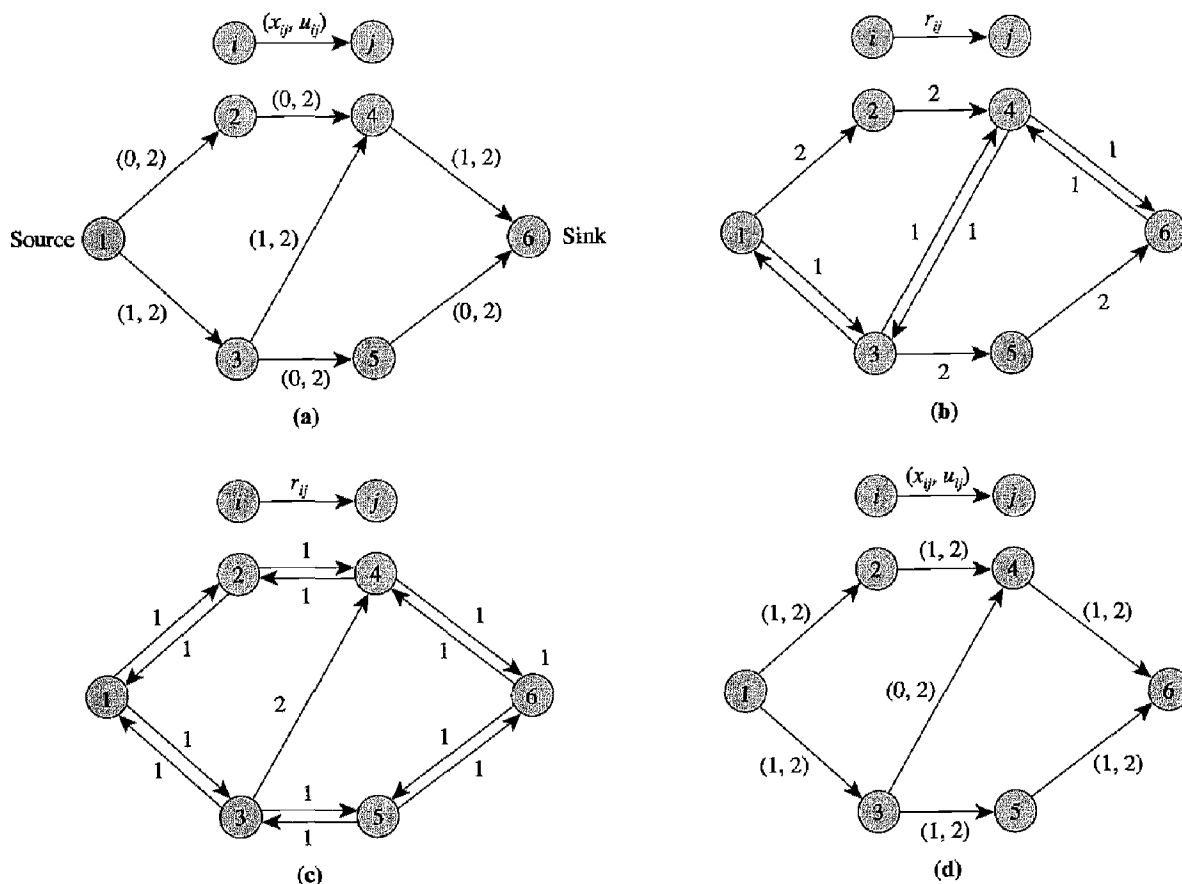


Figure 6.14 The effect of augmentation on flow decomposition: (a) original network with a flow x ; (b) residual network for flow x ; (c) residual network after augmenting one unit along the path 1–2–4–3–5–6; (d) flow in the original network after the augmentation.

one such choice as follows. To highlight this choice, let us rewrite $r_{ij} = u_{ij} - x_{ij} + x_{ji}$ as $x_{ij} - x_{ji} = u_{ij} - r_{ij}$. Now, if $u_{ij} \geq r_{ij}$, we set $x_{ij} = u_{ij} - r_{ij}$ and $x_{ji} = 0$; otherwise, we set $x_{ij} = 0$ and $x_{ji} = r_{ij} - u_{ij}$.

Effect of Augmentation on Flow Decomposition

To obtain better insight concerning the augmenting path algorithm, let us illustrate the effect of an augmentation on the flow decomposition on the preceding example. Figure 6.15(a) gives the decomposition of the initial flow and Figure 6.15(b) gives the decomposition of the flow after we have augmented 1 unit of flow on the path 1–2–4–3–5–6. Although we augmented 1 unit of flow along the path 1–2–4–3–5–6, the flow decomposition contains no such path. Why?

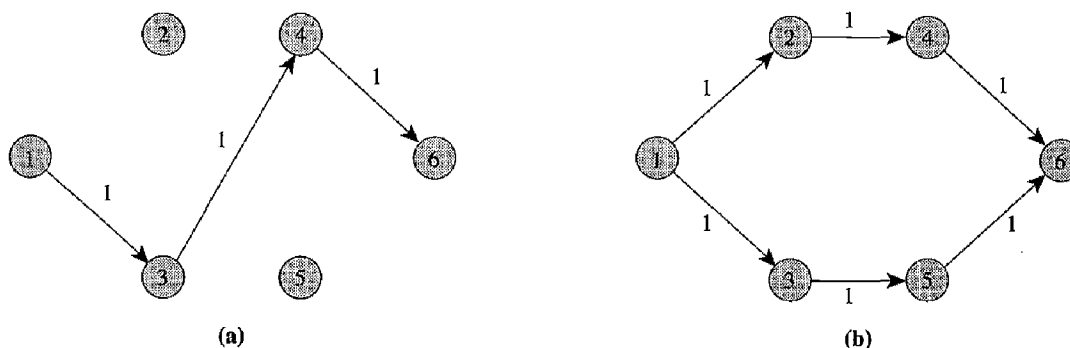


Figure 6.15 Flow decomposition of the solution in (a) Figure 6.14(a) and (b) Figure 6.14(d).

The path 1–3–4–6 defining the flow in Figure 6.14(a) contains three segments: the path up to node 3, arc (3, 4) as a forward arc, and the path up to node 6. We can view this path as an augmentation on the zero flow. Similarly, the path 1–2–4–3–5–6 contains three segments: the path up to node 4, arc (3, 4) as a backward arc, and the path up to node 6. We can view the augmentation on the path 1–2–4–3–5–6 as linking the initial segment of the path 1–3–4–6 with the last segment of the augmentation, linking the last segment of the path 1–3–4–6 with the initial segment of the augmentation, and canceling the flow on arc (3, 4), which then drops from both the path 1–3–4–6 and the augmentation (see Figure 6.16). In general, we can

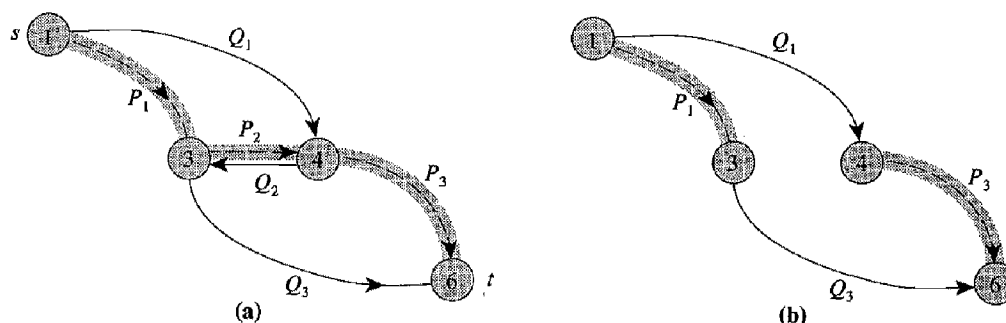


Figure 6.16 The effect of augmentation on flow decomposition: (a) the two augmentations P_1 – P_2 – P_3 and Q_1 – Q_2 – Q_3 ; (b) net effect of these augmentations.

view each augmentation as “pasting together” segments of the current flow decomposition to obtain a new flow decomposition.

6.5 LABELING ALGORITHM AND THE MAX-FLOW MIN-CUT THEOREM

In this section we discuss the augmenting path algorithm in more detail. In our discussion of this algorithm in the preceding section, we did not discuss some important details, such as (1) how to identify an augmenting path or show that the network contains no such path, and (2) whether the algorithm terminates in finite number of iterations, and when it terminates, whether it has obtained a maximum flow. In this section we consider these issues for a specific implementation of the generic augmenting path algorithm known as the *labeling algorithm*. The labeling algorithm is not a polynomial-time algorithm. In Chapter 7, building on the ideas established in this chapter, we describe two polynomial-time implementations of this algorithm.

The labeling algorithm uses a search technique (as described in Section 3.4) to identify a directed path in $G(x)$ from the source to the sink. The algorithm *fans out* from the source node to find all nodes that are reachable from the source along a directed path in the residual network. At any step the algorithm has partitioned the nodes in the network into two groups: *labeled* and *unlabeled*. Labeled nodes are those nodes that the algorithm has reached in the fanning out process and so the algorithm has determined a directed path from the source to these nodes in the residual network; the unlabeled nodes are those nodes that the algorithm has not reached as yet by the fanning-out process. The algorithm iteratively selects a labeled node and scans its arc adjacency list (in the residual network) to reach and label additional nodes. Eventually, the sink becomes labeled and the algorithm sends the maximum possible flow on the path from node s to node t . It then erases the labels and repeats this process. The algorithm terminates when it has scanned all the labeled nodes and the sink remains unlabeled, implying that the source node is not connected to the sink node in the residual network. Figure 6.17 gives an algorithmic description of the labeling algorithm.

Correctness of the Labeling Algorithm and Related Results

To study the correctness of the labeling algorithm, note that in each iteration (i.e., an execution of the whole loop), the algorithm either performs an augmentation or terminates because it cannot label the sink. In the latter case we must show that the current flow x is a maximum flow. Suppose at this stage that S is the set of labeled nodes and $\bar{S} = N - S$ is the set of unlabeled nodes. Clearly, $s \in S$ and $t \in \bar{S}$. Since the algorithm cannot label any node in \bar{S} from any node in S , $r_{ij} = 0$ for each $(i, j) \in (S, \bar{S})$. Furthermore, since $r_{ij} = (u_{ij} - x_{ij}) + x_{ji}$, $x_{ij} \leq u_{ij}$ and $x_{ji} \geq 0$, the condition $r_{ij} = 0$ implies that $x_{ij} = u_{ij}$ for every arc $(i, j) \in (S, \bar{S})$ and $x_{ij} = 0$ for every arc $(i, j) \in (\bar{S}, S)$. [Recall our assumption that for each arc $(i, j) \in A$,

```

algorithm labeling;
begin
    label node  $t$ ;
    while  $t$  is labeled do
        begin
            unlabel all nodes;
            set  $\text{pred}(j) := 0$  for each  $j \in N$ ;
            label node  $s$  and set  $\text{LIST} := \{s\}$ ;
            while  $\text{LIST} \neq \emptyset$  or  $t$  is unlabeled do
                begin
                    remove a node  $i$  from  $\text{LIST}$ ;
                    for each arc  $(i, j)$  in the residual network emanating from node  $i$  do
                        if  $r_{ij} > 0$  and node  $j$  is unlabeled then set  $\text{pred}(j) := i$ , label node  $j$ , and
                            add  $j$  to  $\text{LIST}$ ;
                    end;
                if  $t$  is labeled then augment
            end;
        end;

procedure augment;
begin
    use the predecessor labels to trace back from the sink to the source to
        obtain an augmenting path  $P$  from node  $s$  to node  $t$ ;
     $\delta := \min\{r_{ij} : (i, j) \in P\}$ ;
    augment  $\delta$  units of flow along  $P$  and update the residual capacities;
end;

```

Figure 6.17 Labeling algorithm.

$(j, i) \in A$.] Substituting these flow values in (6.3), we find that

$$v = \sum_{(i,j) \in (S, \bar{S})} x_{ij} - \sum_{(i,j) \in (S, S)} x_{ij} = \sum_{(i,j) \in (S, \bar{S})} u_{ij} = u[S, \bar{S}].$$

This discussion shows that the value of the current flow x equals the capacity of the cut $[S, \bar{S}]$. But then Property 6.1 implies that x is a maximum flow and $[S, \bar{S}]$ is a minimum cut. This conclusion establishes the correctness of the labeling algorithm and, as a by-product, proves the following max-flow min-cut theorem.

Theorem 6.3 (Max-Flow Min-Cut Theorem). *The maximum value of the flow from a source node s to a sink node t in a capacitated network equals the minimum capacity among all s – t cuts.* ♦

The proof of the max-flow min-cut theorem shows that when the labeling algorithm terminates, it has also discovered a minimum cut. The labeling algorithm also proves the following augmenting path theorem.

Theorem 6.4 (Augmenting Path Theorem). *A flow x^* is a maximum flow if and only if the residual network $G(x^*)$ contains no augmenting path.*

Proof. If the residual network $G(x^*)$ contains an augmenting path, clearly the flow x^* is not a maximum flow. Conversely, if the residual network $G(x^*)$ contains no augmenting path, the set of nodes S labeled by the labeling algorithm defines an

s - t cut $[S, \bar{S}]$ whose capacity equals the flow value, thereby implying that the flow must be maximum. ♦

The labeling algorithm establishes one more important result.

Theorem 6.5 (Integrality Theorem). *If all arc capacities are integer, the maximum flow problem has an integer maximum flow.*

Proof. This result follows from an induction argument applied to the number of augmentations. Since the labeling algorithm starts with a zero flow and all arc capacities are integer, the initial residual capacities are all integer. The flow augmented in any iteration equals the minimum residual capacity of some path, which by the induction hypothesis is integer. Consequently, the residual capacities in the next iteration will again be integer. Since the residual capacities r_{ij} and the arc capacities u_{ij} are all integer, when we convert the residual capacities into flows by the method described previously, the arc flows x_{ij} will be integer valued as well. Since the capacities are integer, each augmentation adds at least one unit to the flow value. Since the maximum flow cannot exceed the capacity of any cut, the algorithm will terminate in a finite number of iterations. ♦

The integrality theorem does not imply that every optimal solution of the maximum flow problem is integer. The maximum flow problem may have noninteger solutions and, most often, has such solutions. The integrality theorem shows that the problem always has at least one integer optimal solution.

Complexity of the Labeling Algorithm

To study the worst-case complexity of the labeling algorithm, recall that in each iteration, except the last, when the sink cannot be labeled, the algorithm performs an augmentation. It is easy to see that each augmentation requires $O(m)$ time because the search method examines any arc or any node at most once. Therefore, the complexity of the labeling algorithm is $O(m)$ times the number of augmentations. How many augmentations can the algorithm perform? If all arc capacities are integral and bounded by a finite number U , the capacity of the cut $(s, N - \{s\})$ is at most nU . Therefore, the maximum flow value is bounded by nU . The labeling algorithm increases the value of the flow by at least 1 unit in any augmentation. Consequently, it will terminate within nU augmentations, so $O(nmU)$ is a bound on the running time of the labeling algorithm. Let us formally record this observation.

Theorem 6.6. *The labeling algorithm solves the maximum flow problem in $O(nmU)$ time.* ♦

Throughout this section, we have assumed that each arc capacity is finite. In some applications, it will be convenient to model problems with infinite capacities on some arcs. If we assume that some s - t cut has a finite capacity and let U denote the maximum capacity across this cut, Theorem 6.6 and, indeed, all the other results in this section remain valid. Another approach for addressing situations with infinite capacity arcs would be to impose a capacity on these arcs, chosen sufficiently large

as to not affect the maximum flow value (see Exercise 6.23). In defining the residual capacities and developing algorithms to handle situations with infinite arc capacities, we adopt this approach rather than modifying the definitions of residual capacities.

Drawbacks of the Labeling Algorithm

The labeling algorithm is possibly the simplest algorithm for solving the maximum flow problem. Empirically, the algorithm performs reasonably well. However, the worst-case bound on the number of iterations is not entirely satisfactory for large values of U . For example, if $U = 2^n$, the bound is exponential in the number of nodes. Moreover, the algorithm can indeed perform this many iterations, as the example given in Figure 6.18 illustrates. For this example, the algorithm can select the augmenting paths $s-a-b-t$ and $s-b-a-t$ alternatively 10^6 times, each time augmenting unit flow along the path. This example illustrates one shortcoming of the algorithm.

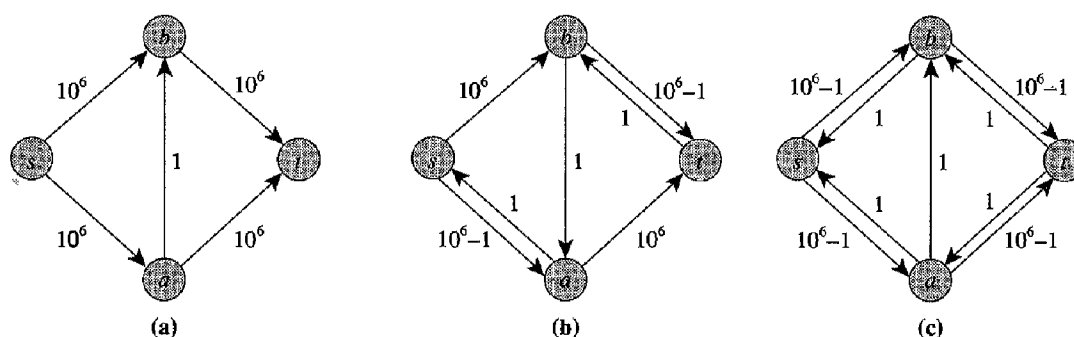


Figure 6.18 Pathological example of the labeling algorithm: (a) residual network for the zero flow; (b) network after augmenting unit flow along the path $s-a-b-t$; (c) network after augmenting unit flow along the path $s-b-a-t$.

A second drawback of the labeling algorithm is that if the capacities are irrational, the algorithm might not terminate. For some pathological instances of the maximum flow problem (see Exercise 6.48), the labeling algorithm does not terminate, and although the successive flow values converge, they converge to a value strictly less than the maximum flow value. (Note, however, that the max-flow min-cut theorem holds even if arc capacities are irrational.) Therefore, if the labeling algorithm is guaranteed to be effective, it must select augmenting paths carefully.

A third drawback of the labeling algorithm is its “forgetfulness.” In each iteration, the algorithm generates node labels that contain information about augmenting paths from the source to other nodes. The implementation we have described erases the labels as it moves from one iteration to the next, even though much of this information might be valid in the next iteration. Erasing the labels therefore destroys potentially useful information. Ideally, we should retain a label when we can use it profitably in later computations.

In Chapter 7 we describe several improvements of the labeling algorithm that overcomes some or all of these drawbacks. Before discussing these improvements, we discuss some interesting implications of the max-flow min-cut theorem.

6.6 COMBINATORIAL IMPLICATIONS OF THE MAX-FLOW MIN-CUT THEOREM

As we noted in Section 6.2 when we discussed several applications of the maximum flow problem, in some applications we wish to find a minimum cut in a network, which we now know is equivalent to finding a maximum flow in the network. In fact, the relationship between maximum flows and minimum cuts permits us to view many problems from either of two dual perspectives: a flow perspective or a cut perspective. At times this dual perspective provides novel insight about an underlying problem. In particular, when applied in various ways, the max-flow min-cut theorem reduces to a number of min-max duality relationships in combinatorial theory. In this section we illustrate this use of network flow theory by developing several results in combinatorics. We might note that these results are fairly deep and demonstrate the power of the max-flow min-cut theorem. To appreciate the power of the max-flow min-cut theorem, we would encourage the reader to try to prove the following results without using network flow theory.

Network Connectivity

We first study some connectivity issues about networks that arise, for example, in the design of communication networks. We first define some notation. We refer to two directed paths from node s to node t as *arc disjoint* if they do not have any arc in common. Similarly, we refer to two directed paths from node s to node t as *node disjoint* if they do not have any node in common, except the source and the sink nodes. Given a directed network $G = (N, A)$ and two specified nodes s and t , we are interested in the following two questions: (1) What is the maximum number of arc-disjoint (directed) paths from node s to node t ; and (2) what is the minimum number of arcs that we should remove from the network so that it contains no directed paths from node s to node t ? The following theorem shows that these two questions are really alternative ways to address the same issue.

Theorem 6.7. *The maximum number of arc-disjoint paths from node s to node t equals the minimum number of arcs whose removal from the network disconnects all paths from node s to node t .*

Proof. Define the capacity of each arc in the network as equal to 1. Consider any feasible flow x of value v in the resulting unit capacity network. The flow decomposition theorem (Theorem 3.5) implies that we can decompose the flow x into flows along paths and cycles. Since flows around cycles do not affect the flow value, the flows on the paths sum to v . Furthermore, since each arc capacity is 1, these paths are arc disjoint and each carries 1 unit of flow. Consequently, the network contains v arc-disjoint paths from s to t .

Now consider any s - t cut $[S, \bar{S}]$ in the network. Since each arc capacity is 1, the capacity of this cut is $| (S, \bar{S}) |$ (i.e., it equals the number of forward arcs in the cut). Since each path from node s to node t contains at least one arc in (S, \bar{S}) , the removal of the arcs in (S, \bar{S}) disconnects all the paths from node s to node t . Consequently, the network contains a disconnecting set of arcs of cardinality equal

to the capacity of any s - t cut $[S, \bar{S}]$. The max-flow min-cut theorem immediately implies that the maximum number of arc-disjoint paths from s to t equals the minimum number of arcs whose removal will disconnect all paths from node s to node t .

◆

We next discuss the node-disjoint version of the preceding theorem.

Theorem 6.8. *The maximum number of node-disjoint paths from node s to node t equals the minimum number of nodes whose removal from the network disconnects all paths from nodes s to node t .*

Proof. Split each node i in G , other than s and t , into two nodes i' and i'' and add a “node-splitting” arc (i', i'') of unit capacity. All the arcs in G entering node i now enter node i' and all the arcs emanating from node i now emanate from node i'' . Let G' denote this transformed network. Assign a capacity of ∞ to each arc in the network except the node-splitting arcs, which have unit capacity. It is easy to see that there is one-to-one correspondence between the arc-disjoint paths in G' and the node-disjoint paths in G . Therefore, the maximum number of arc-disjoint paths in G' equals the maximum number of node-disjoint paths in G .

As in the proof of Theorem 6.7, flow decomposition implies that a flow of v units from node s to node t in G' decomposes into v arc-disjoint paths each carrying unit flow; and these v arc-disjoint paths in G' correspond to v node-disjoint paths in G . Moreover, note that any s - t cut with finite capacity contains only node-splitting arcs since all other arcs have infinite capacity. Therefore, any s - t cut in G' with capacity k corresponds to a set of k nodes whose removal from G destroys all paths from node s to node t . Applying the max-flow min-cut theorem to G' and using the preceding observations establishes that the maximum number of node-disjoint paths in G from node s to node t equals the minimum number of nodes whose removal from G disconnects nodes s and t .

◆

Matchings and Covers

We next state some results about matchings and node covers in a bipartite network. For a directed bipartite network $G = (N_1 \cup N_2, A)$ we refer to a subset $A' \subseteq A$ as a *matching* if no two arcs in A' are incident to the same node (i.e., they do not have any common endpoint). We refer to a subset $N' \subseteq N = N_1 \cup N_2$ as a *node cover* if every arc in A is incident to one of the nodes in N' . For illustrations of these definitions, consider the bipartite network shown in Figure 6.19. In this network the set of arcs $\{(1, 1'), (3, 3'), (4, 5'), (5, 2')\}$ is a matching but the set of arcs $\{(1, 2'), (3, 1'), (3, 4')\}$ is not because the arcs $(3, 1')$ and $(3, 4')$ are incident to the same node 3. In the same network the set of nodes $\{1, 2', 3, 5'\}$ is a node cover, but the set of nodes $\{2', 3', 4, 5\}$ is not because the arcs $(1, 1')$, $(3, 1')$, and $(3, 4')$ are not incident to any node in the set.

Theorem 6.9. *In a bipartite network $G = (N_1 \cup N_2, A)$, the maximum cardinality of any matching equals the minimum cardinality of any node cover of G .*

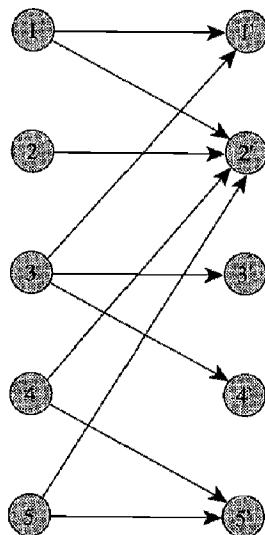


Figure 6.19 Bipartite network.

Proof. Augment the network by adding a source node s and an arc (s, i) of unit capacity for each $i \in N_1$. Similarly, add a sink node t and an arc (j, t) of unit capacity for each $j \in N_2$. Denote the resulting network by G' . We refer to the arcs in A as *original arcs* and the additional arcs as *artificial arcs*. We set the capacity of each artificial arc equal to 1 and the capacity of each original arc equal to ∞ .

Now consider any flow x of value v from node s to node t in the network G' . We can decompose the flow x into v paths of the form $s-i-j-t$ each carrying 1 unit of flow. Thus v arcs of the original network have a positive flow. Furthermore, these arcs constitute a matching, for otherwise the flow on some artificial arc would exceed 1 unit. Consequently, a flow of value v corresponds to a matching of cardinality v . Similarly, a matching of cardinality v defines a flow of value v .

We next show that any node cover H of $G = (N_1 \cup N_2, A)$ defines an $s-t$ cut of capacity $|H|$ in G' . Given the node cover H , construct a set of arcs Q as follows: For each $i \in H$, if $i \in N_1$, add arc (s, i) to Q , and if $i \in N_2$, add arc (i, t) to Q . Since H is a node cover, each directed path from node s to node t in G' contains one arc in Q ; therefore, Q is a valid $s-t$ cut of capacity $|H|$.

We now show the converse result; that is, for a given $s-t$ cut Q of capacity k in G' , the network G contains a node cover of cardinality k . We first note that the cut Q consists solely of artificial arcs because the original arcs have infinite capacity. From Q we construct a set H of nodes as follows: if $(s, i) \in Q$ and $(i, t) \in Q$, we add i to H . Now observe that each original arc (i, j) defines a directed path $s-i-j-t$ in G' . Since Q is an $s-t$ cut, either $(s, i) \in Q$ or $(j, t) \in Q$ or both. By the preceding construction, either $i \in H$ or $j \in H$ or both. Consequently, H must be a node cover. We have thus established a one-to-one correspondence between node covers in G and $s-t$ cuts in G' .

The max-flow min-cut theorem implies that the maximum flow value equals the capacity of a minimum cut. In view of the max-flow min-cut theorem, the preceding observations imply that the maximum number of independent arcs in G equals the minimum number of nodes in a node cover of G . The theorem thus follows. ♦

Figure 6.20 gives a further illustration of Theorem 6.9. In this figure, we have transformed the matching problem of Figure 6.19 into a maximum flow problem, and we have identified the minimum cut. The minimum cut consists of the arcs $(s, 1)$, $(s, 3)$, $(2', t)$ and $(5', t)$. Correspondingly, the set $\{1, 3, 2', 5'\}$ is a minimum cardinality node cover, and a maximum cardinality matching is $(1, 1')$, $(2, 2')$, $(3, 3')$ and $(5, 5')$.

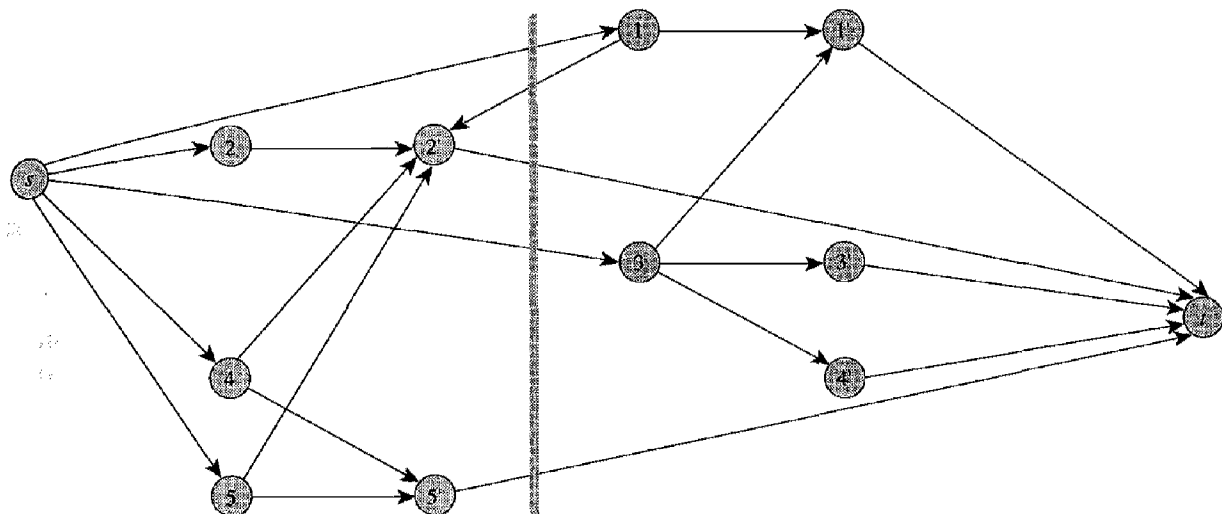


Figure 6.20 Minimum cut for the maximum flow problem defined in Figure 6.19.

As we have seen in the discussion throughout this section, the max-flow min-cut theorem is a powerful tool for establishing a number of results in the field of combinatorics. Indeed, the range of applications of the max-flow min-cut theorem and the ability of this theorem to encapsulate so many subtle duality (i.e., max-min) results as special cases is quite surprising, given the simplicity of the labeling algorithm and of the proof of the max-flow min-cut theorem. The wide range of applications reflects the fact that flows and cuts, and the relationship between them, embody central combinatorial results in many problem domains within applied mathematics.

6.7 FLOWS WITH LOWER BOUNDS

In this section we consider maximum flow problems with nonnegative lower bounds imposed on the arc flows; that is, the flow on any arc $(i, j) \in A$ must be at least $l_{ij} \geq 0$. The following formulation models this problem:

Maximize v

subject to

$$\sum_{\{j: (i,j) \in A\}} x_{ij} - \sum_{\{j: (j,i) \in A\}} x_{ji} = \begin{cases} v & \text{for } i = s, \\ 0 & \text{for all } i \in N - \{s, t\}, \\ -v & \text{for } i = t, \end{cases}$$

$$l_{ij} \leq x_{ij} \leq u_{ij} \quad \text{for each } (i, j) \in A.$$

In previous sections we studied a special case of this problem with only zero lower bounds. Whereas the maximum flow problem with zero lower bounds always has a feasible solution (since the zero flow is feasible), the problem with nonnegative lower bounds could be infeasible. For example, consider the maximum flow problem given in Figure 6.21. This problem does not have a feasible solution because arc $(1, 2)$ must carry at least 5 units of flow into node 2 and arc $(2, 3)$ can remove at most 4 units of flow; therefore, we can never satisfy the mass balance constraint of node 2.

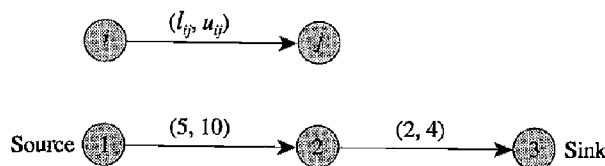


Figure 6.21 Maximum flow problem with no feasible solution.

As illustrated by this example, any maximum flow algorithm for problems with nonnegative lower bounds has two objectives: (1) to determine whether the problem is feasible, and (2) if so, to establish a maximum flow. It therefore comes as no surprise that most algorithms use a two-phase approach. The first phase determines a feasible flow if one exists, and the second phase converts a feasible flow into a maximum flow. We shall soon see that the problem in each phase essentially reduces to solving a maximum flow problem with zero lower bounds. Consequently, it is possible to solve the maximum flow problem with nonnegative lower bounds by solving two maximum flow problems, each with zero lower bounds. For convenience, we consider the second phase prior to the first phase.

Determining a Maximum Flow

Suppose that we have a feasible flow x in the network. We can then modify any maximum flow algorithm designed for the zero lower bound case to obtain a maximum flow. In these algorithms, we make only one modification: We define the residual capacity of an arc (i, j) as $r_{ij} = (u_{ij} - x_{ij}) + (x_{ji} - l_{ji})$; the first term in this expression denotes the maximum increase in flow from node i to node j using the remaining capacity of arc (i, j) , and the second term denotes the maximum increase in flow from node i to node j by canceling the existing flow on arc (j, i) . Notice that since each arc flow is within its lower and upper bounds, each residual capacity is nonnegative. Recall that the maximum flow algorithm described in this chapter (and the ones described in Chapter 7) works with only residual capacities and does not need arc flows, capacities, or lower bounds. Therefore we can use any of these algorithms to establish a maximum flow in the network. These algorithms terminate with optimal residual capacities. From these residual capacities we can construct maximum flow in a large number of ways. For example, through a change of variables we can reduce the computations to a situation we have considered before. For all arcs (i, j) , let $u'_{ij} = u_{ij} - l_{ij}$, $r'_{ij} = r_{ij}$, and $x'_{ij} = x_{ij} - l_{ij}$. The residual capacity for arc (i, j) is $r_{ij} = (u_{ij} - x_{ij}) + (x_{ji} - l_{ji})$. Equivalently, $r'_{ij} = u'_{ij} - x'_{ij} + x'_{ji}$. Similarly, $r'_{ji} = u'_{ji} - x'_{ji} + x'_{ij}$. If we compute the x' values in terms of r' and u' , we obtain the same expression as before, i.e., $x'_{ij} = \max(u'_{ij} - r'_{ij}, 0)$ and $x'_{ji} = \max(u'_{ji} -$

$r'_{ij}, 0$). Converting back into the original variables, we obtain the following formulae:

$$\begin{aligned}x_{ij} &= l_{ij} + \max(u_{ij} - r_{ij} - l_{ij}, 0), \\x_{ji} &= l_{ji} + \max(u_{ji} - r_{ji} - l_{ji}, 0).\end{aligned}$$

We now show that the solution determined by this modified procedure solves the maximum flow problem with nonnegative lower bounds. Let x denote a feasible flow in G with value equal to v . Moreover, let $[S, \bar{S}]$ denote an s - t cut. We define the *capacity* of an s - t cut $[S, \bar{S}]$ as

$$u[S, \bar{S}] = \sum_{(i,j) \in (S, \bar{S})} u_{ij} - \sum_{(i,j) \in (\bar{S}, S)} l_{ij}. \quad (6.7)$$

The capacity of the cut denotes the maximum amount of “net” flow that can be sent out of the node set S . We next use equality (6.3), which we restate for convenience.

$$v = \sum_{(i,j) \in (S, \bar{S})} x_{ij} - \sum_{(i,j) \in (\bar{S}, S)} x_{ij}. \quad (6.8)$$

Substituting $x_{ij} \leq u_{ij}$ in the first summation and $l_{ij} \leq x_{ij}$ in the second summation of this inequality shows that

$$v \leq \sum_{(i,j) \in (S, \bar{S})} u_{ij} - \sum_{(i,j) \in (\bar{S}, S)} l_{ij} = u[S, \bar{S}]. \quad (6.9)$$

Inequality (6.9) indicates that the maximum flow value is less than or equal to the capacity of any s - t cut. At termination, the maximum flow algorithm obtains an s - t cut $[S, \bar{S}]$ with $r_{ij} = 0$ for every arc $(i, j) \in (S, \bar{S})$. Let x denote the corresponding flow with value equal to v . Since $r_{ij} = (u_{ij} - x_{ij}) + (x_{ji} - l_{ji})$, the conditions $x_{ij} \leq u_{ij}$ and $l_{ji} \leq x_{ji}$ imply that $x_{ij} = u_{ij}$ and $x_{ji} = l_{ji}$. Consequently, $x_{ij} = u_{ij}$ for every arc $(i, j) \in (S, \bar{S})$ and $x_{ij} = l_{ij}$ for every arc $(i, j) \in (\bar{S}, S)$. Substituting these values in (6.8), we find that

$$v = u[S, \bar{S}] = \sum_{(i,j) \in (S, \bar{S})} u_{ij} - \sum_{(i,j) \in (\bar{S}, S)} l_{ij}. \quad (6.10)$$

In view of inequality (6.9), equation (6.10) implies that $[S, \bar{S}]$ is a minimum s - t cut and x is a maximum flow. As a by-product of this result, we have proved a generalization of the max-flow min-cut theorem for problems with nonnegative lower bounds.

Theorem 6.10 (Generalized Max-Flow Min-Cut Theorem). *If the capacity of an s - t cut $[S, \bar{S}]$ in a network with both lower and upper bounds on arc flows is defined by (6.7), the maximum value of flow from node s to node t equals the minimum capacity among all s - t cuts.* ♦

Establishing a Feasible Flow

We now address the issue of determining a feasible flow in the network. We first transform the maximum flow problem into a circulation problem by adding an arc (t, s) of infinite capacity. This arc carries the flow sent from node s to node t back to node s . Consequently, in the circulation formulation of the problem, the outflow

of each node, including nodes s and t , equals its flow. Clearly, the maximum flow problem admits a feasible flow if and only if the circulation problem admits a feasible flow. Given the possibility of making this transformation, we now focus our intention on finding a feasible circulation, and characterizing conditions when an arbitrary circulation problem, with lower and upper bounds of flows, possesses a feasible solution.

The feasible circulation problem is to identify a flow x satisfying the following constraints:

$$\sum_{\{j:(i,j) \in A\}} x_{ij} - \sum_{\{j:(j,i) \in A\}} x_{ji} = 0 \quad \text{for all } i \in N, \quad (6.11a)$$

$$l_{ij} \leq x_{ij} \leq u_{ij} \quad \text{for all } (i, j) \in A. \quad (6.11b)$$

By replacing $x_{ij} = x'_{ij} + l_{ij}$ in constraints (6.11a) and (6.11b), we obtain the following transformed problem:

$$\sum_{\{j:(i,j) \in A\}} x'_{ij} - \sum_{\{j:(j,i) \in A\}} x'_{ji} = b(i) \quad \text{for all } i \in N, \quad (6.12a)$$

$$0 \leq x'_{ij} \leq u_{ij} - l_{ij} \quad \text{for all } (i, j) \in A, \quad (6.12b)$$

with supplies/demands $b(\cdot)$ at the nodes defined by

$$b(i) = \sum_{\{j:(j,i) \in A\}} l_{ji} - \sum_{\{j:(i,j) \in A\}} l_{ij}.$$

Observe that $\sum_{i \in N} b(i) = 0$ since each l_{ij} occurs twice in this expression, once with a positive sign and once with a negative sign. The feasible circulation problem is then equivalent to determining whether the transformed problem has a solution x' satisfying (6.12).

Notice that this problem is essentially the same as the feasible flow problem discussed in Application 6.1. In discussing this application we showed that by solving a maximum flow problem we either determine a solution satisfying (6.12) or show that no solution satisfies (6.12). If x'_{ij} is a feasible solution of (6.12), $x_{ij} = x'_{ij} + l_{ij}$ is a feasible solution of (6.11).

Characterizing a Feasible Flow

We next characterize feasible circulation problems (i.e., derive the necessary and sufficiency conditions for a circulation problem to possess a feasible solution). Let S be any set of nodes in the network. By summing the mass balance constraints of the nodes in S , we obtain the expression

$$\sum_{(i,j) \in (S, \bar{S})} x_{ij} - \sum_{(i,j) \in (\bar{S}, S)} x_{ij} = 0. \quad (6.13)$$

Using the inequalities $x_{ij} \leq u_{ij}$ in the first term of (6.13) and the inequalities $x_{ij} \geq l_{ij}$ in the second term, we find that

$$\sum_{(i,j) \in (S, \bar{S})} l_{ij} \leq \sum_{(i,j) \in (\bar{S}, S)} u_{ij}. \quad (6.14)$$

The expression in (6.14), which is a necessary condition for feasibility, states that

the maximum amount of flow that we can send out from a set S of nodes must be at least as large as the minimum amount of flow that the nodes in S must receive. Clearly, if a set of nodes must receive more than what the other nodes can send them, the network has no feasible circulation. As we will see, these conditions are also sufficient for ensuring feasibility [i.e., if the network data satisfies the conditions (6.14) for every set S of nodes, the network has a feasible circulation that satisfies the flow bounds on all its arcs].

We give an algorithmic proof for the sufficiency of condition (6.14). The algorithm starts with a circulation x that satisfies the mass balance and capacity constraints, but might violate some of the lower bound constraints. The algorithm gradually converts this circulation into a feasible flow or identifies a node set S that violates condition (6.14).

With respect to a flow x , we refer to an arc (i, j) as *infeasible* if $x_{ij} < l_{ij}$ and *feasible* if $l_{ij} \leq x_{ij}$. The algorithm selects an infeasible arc (p, q) and attempts to make it feasible by increasing the flow on this arc. The mass balance constraints imply that in order to increase the flow on the arc, we must augment flow along one or more cycles in the residual network that contain arc (p, q) as a forward arc. We define the residual network $G(x)$ with respect to a flow x the same way we defined it previously except that we set the residual capacity of any infeasible arc (i, j) to the value $u_{ij} - x_{ij}$. Any augmenting cycle containing arc (p, q) as a forward arc must consist of a directed path in the residual network $G(x)$ from node q to node p plus the arc (p, q) . We can use a labeling algorithm to identify a directed path from node q to node p .

We apply this procedure to one infeasible arc at a time, at each step decreasing the infeasibility of the arcs until we either identify a feasible flow or the labeling algorithm is unable to identify a directed path from node q to node p for some infeasible arc (p, q) . We show that in the latter case, the maximum flow problem must be infeasible. Let S be the set of nodes labeled by the last application of the labeling algorithm. Clearly, $q \in S$ and $p \in \bar{S} \equiv N - S$. Since the labeling algorithm cannot label any node not in S , every arc (i, j) from S to \bar{S} has a residual capacity of value zero. Therefore, $x_{ij} = u_{ij}$ for every arc $(i, j) \in (S, \bar{S})$ and $x_{ij} \leq l_{ij}$ for every arc $(i, j) \in (\bar{S}, S)$. Also observe that $(p, q) \in (\bar{S}, S)$ and $x_{pq} < l_{pq}$. Substituting these values in (6.13), we find that

$$\sum_{(i,j) \in (\bar{S}, S)} l_{ij} > \sum_{(i,j) \in (S, \bar{S})} u_{ij},$$

contradicting condition (6.14), which we have already shown is necessary for feasibility. We have thus established the following fundamental theorem.

Theorem 6.11 (Circulation Feasibility Conditions). *A circulation problem with nonnegative lower bounds on arc flows is feasible if and only if for every set S of nodes*

$$\sum_{(i,j) \in (\bar{S}, S)} l_{ij} \leq \sum_{(i,j) \in (S, \bar{S})} u_{ij}. \quad \blacklozenge$$

Note that the proof of this theorem specifies a one pass algorithm, starting with the zero flow, for finding a feasible solution to any circulation problem whose arc

upper bounds u_{ij} are all nonnegative. In Exercise 6.7 we ask the reader to specify a one pass algorithm for any situation (i.e., even when some upper bounds are negative).

A by-product of Theorem 6.11 is the following result, which states necessary and sufficiency conditions for the existence of a feasible solution for the feasible flow problem stated in (6.2). (Recall that a feasible flow problem is the feasibility version of the minimum cost flow problem.) We discuss the proof of this result in Exercise 6.43.

Theorem 6.12. *The feasible flow problem stated in (6.2) has a feasible solution if and only if for every subset $S \subseteq N$, $b(S) - u[S, \bar{S}] \leq 0$, where $b(S) = \sum_{i \in S} b(i)$.* ♦

6.8 SUMMARY

In this chapter we studied two closely related problems: the maximum flow problem and the minimum cut problem. After illustrating a variety of applications of these problems, we showed that the maximum flow and the minimum cut problems are closely related of each other (in fact, they are dual problems) and solving the maximum flow problem also solves the minimum cut problem. We began by showing that the value of any flow is less than or equal to the capacity of any cut in the network (i.e., this is a “weak duality” result). The fact that the value of some flow equals the capacity of some cut in the network (i.e., the “strong duality” result) is a deeper result. This result is known as the *max-flow min-cut theorem*. We establish it by specifying a labeling algorithm that maintains a feasible flow x in the network and sends additional flow along directed paths from the source node to the sink node in the residual network $G(x)$. Eventually, $G(x)$ contains no directed path from the source to the sink. At this point, the value of the flow x equals the capacity of some cut $[S, \bar{S}]$ in the network. The weak duality result implies that x is a maximum flow and $[S, \bar{S}]$ is a minimum cut. Since the labeling algorithm maintains an integer flow at every step (assuming integral capacity data), the optimal flow that it finds is integral. This result is a special case of a more general network flow integrality result that we establish in Chapter 9. The labeling algorithm runs in $O(nmU)$ time. This time bound is not attractive from the worst-case perspective. In Chapter 7 we develop two polynomial-time implementations of the labeling algorithm.

The max-flow min-cut theorem has far-reaching implications. It allows us to prove several important results in combinatorics that appear difficult to prove using other means. We proved the following results: (1) the maximum number of arc-disjoint (or node-disjoint) paths connecting two nodes s and t in a network equals the minimum number of arcs (or nodes) whose removal from the network leaves no directed path from node s to node t ; and (2) in a bipartite network, the maximum cardinality of any matching equals the minimum cardinality of any node cover. In the exercises we ask the reader to prove other implications of the max-flow min-cut theorem.

To conclude this chapter we studied the maximum flow problem with non-negative lower bounds on arc flows. We can solve this problem using a two-phase approach. The first phase determines a feasible flow if one exists, and the second phase converts this flow into a maximum flow; in both phases we solve a maximum

flow problem with zero lower bounds. We also described a theoretical result for characterizing when a maximum flow problem with nonnegative lower bounds has a feasible solution. Roughly speaking, this characterization states that the maximum flow problem has a feasible solution if and only if the maximum possible outflow of every cut is at least as large as the minimum required inflow for that cut.

REFERENCE NOTES

The seminal paper of Ford and Fulkerson [1956a] on the maximum flow problem established the celebrated max-flow min-cut theorem. Fulkerson and Dantzig [1955], and Elias, Feinstein, and Shannon [1956] independently established this result. Ford and Fulkerson [1956a] and Elias et al. [1956] solved the maximum flow problem by augmenting path algorithms, whereas Fulkerson and Dantzig [1955] solved it by specializing the simplex method for linear programming. The labeling algorithm that we described in Section 6.5 is due to Ford and Fulkerson [1956a]; their classical book, Ford and Fulkerson [1962], offers an extensive treatment of this algorithm. Unfortunately, the labeling algorithm runs in pseudopolynomial time; moreover, as shown by Ford and Fulkerson [1956a], for networks with arbitrary irrational arc capacities, the algorithm can perform an infinite sequence of augmentations and might converge to a value different from the maximum flow value. Several improved versions of the labeling algorithm overcome this limitation. We provide citations to these algorithms and to their improvements in the reference notes of Chapter 7. In Chapter 7 we also discuss computational properties of maximum flow algorithms.

In Section 6.6 we studied the combinatorial implications of the max-flow min-cut theorem. Theorems 6.7 and 6.8 are known as Menger's theorem. Theorem 6.9 is known as the König-Egerváry theorem. Ford and Fulkerson [1962] discuss these and several additional combinatorial results that are provable using the max-flow min-cut theorem.

In Section 6.7 we studied the feasibility of a network flow problem with non-negative lower bounds imposed on the arc flows. Theorem 6.11 is due to Hoffman [1960], and Theorem 6.12 is due to Gale [1957]. The book by Ford and Fulkerson [1962] discusses these and some additional feasibility results extensively. The algorithm we have presented for identifying a feasible flow in a network with non-negative lower bounds is adapted from this book.

The applications of the maximum flow problem that we described in Section 6.2 are adapted from the following papers:

1. Feasible flow problem (Berge and Ghouila-Houri [1962])
2. Problem of representatives (Hall [1956])
3. Matrix rounding problem (Bacharach [1966])
4. Scheduling on uniform parallel machines (Federgruen and Groenevelt [1986])
5. Distributed computing on a two-processor model (Stone [1977])
6. Tanker scheduling problem (Dantzig and Fulkerson [1954])

Elsewhere in this book we describe other applications of the maximum flow problem. These applications include: (1) the tournament problem (Application 1.3, Ford and Johnson [1959]), (2) the police patrol problem (Exercise 1.9, Khan [1979]),

(3) nurse staff scheduling (Exercise 6.2, Khan and Lewis [1987]), (4) solving a system of equations (Exercise 6.4, Lin [1986]), (5) statistical security of data (Exercises 6.5, Application 8.3, Gusfield [1988], Kelly, Golden, and Assad [1990]), (6) the minimax transportation problem (Exercise 6.6, Ahuja [1986]), (7) the baseball elimination problem (Application 8.1, Schwartz [1966]), (8) network reliability testing (Application 8.2, Van Slyke and Frank [1972]), (9) open pit mining (Application 19.1, Johnson [1968]), (10) selecting freight handling terminals (Application 19.2, Rhys [1970]), (11) optimal destruction of military targets (Application 19.3, Orlin [1987]), (12) the flyaway kit problem (Application 19.4, Mamer and Smith [1982]), (13) maximum dynamic flows (Application 19.12, Ford and Fulkerson [1958a]), and (14) models for building evacuation (Application 19.13, Chalmet, Francis, and Saunders [1982]).

Two other interesting applications of the maximum flow problem are preemptive scheduling on machines with different speeds (Martel [1982]), and the multifacility rectilinear distance location problem (Picard and Ratliff [1978]). The following papers describe additional applications or provide additional references: McGinnis and Nuttle [1978], Picard and Queyranne [1982], Abdallaoui [1987], Gusfield, Martel, and Fernandez-Baca [1987], Gusfield and Martel [1989], and Gallo, Grigoriadis, and Tarjan [1989].

EXERCISES

- 6.1. Dining problem.** Several families go out to dinner together. To increase their social interaction, they would like to sit at tables so that no two members of the same family are at the same table. Show how to formulate finding a seating arrangement that meets this objective as a maximum flow problem. Assume that the dinner contingent has p families and that the i th family has $a(i)$ members. Also assume that q tables are available and that the j th table has a seating capacity of $b(j)$.
- 6.2. Nurse staff scheduling** (Khan and Lewis [1987]). To provide adequate medical service to its constituents at a reasonable cost, hospital administrators must constantly seek ways to hold staff levels as low as possible while maintaining sufficient staffing to provide satisfactory levels of health care. An urban hospital has three departments: the emergency room (department 1), the neonatal intensive care nursery (department 2), and the orthopedics (department 3). The hospital has three work shifts, each with different levels of necessary staffing for nurses. The hospital would like to identify the minimum number of nurses required to meet the following three constraints: (1) the hospital must allocate at least 13, 32, and 22 nurses to the three departments (over all shifts); (2) the hospital must assign at least 26, 24, and 19 nurses to the three shifts (over all departments); and (3) the minimum and maximum number of nurses allocated to each department in a specific shift must satisfy the following limits:

		Department		
		1	2	3
Shift	1	(6, 8)	(11, 12)	(7, 12)
	2	(4, 6)	(11, 12)	(7, 12)
	3	(2, 4)	(10, 12)	(5, 7)

Suggest a method using maximum flows to identify the minimum number of nurses required to satisfy all the constraints.

- 6.3. A commander is located at one node p in a communication network G and his subordinates are located at nodes denoted by the set S . Let u_{ij} be the effort required to eliminate arc (i, j) from the network. The problem is to determine the minimal effort required to block all communications between the commander and his subordinates. How can you solve this problem in polynomial time?
- 6.4. **Solving a system of equations** (Lin [1986]). Let $F = \{f_{ij}\}$ be a given $p \times q$ matrix and consider the following system of $p + q$ equations in the (possibly fractional) variables y :

$$\sum_{j=1}^q f_{ij} y_{ij} = u_i, \quad 1 \leq i \leq p, \quad (6.15a)$$

$$\sum_{i=1}^p f_{ij} y_{ij} = v_j, \quad 1 \leq j \leq q. \quad (6.15b)$$

In this system $u_i \geq 0$ and $v_j \geq 0$ are given constants satisfying the condition $\sum_{i=1}^p u_i = \sum_{j=1}^q v_j$.

- (a) Define a matrix $D = \{d_{ij}\}$ as follows: $d_{ij} = 0$ if $f_{ij} = 0$, and $d_{ij} = 1$ if $f_{ij} \neq 0$. Show that (6.15) has a feasible solution if and only if the following system of $p + q$ equations has a feasible solution x :

$$\sum_{j=1}^q d_{ij} x_{ij} = u_i, \quad 1 \leq i \leq p, \quad (6.16a)$$

$$\sum_{i=1}^p d_{ij} x_{ij} = v_j, \quad 1 \leq j \leq q. \quad (6.16b)$$

- (b) Show how to formulate the problem of identifying a feasible solution of the system (6.16) as a feasible circulation problem (i.e., identifying a circulation in some network with lower and upper bounds imposed on the arc flows). [Hint: The network has a node i for the i th row in (6.16a), a node j for the j th row in (6.16b), and one extra node s .]
- 6.5. **Statistical security of data** (Kelly, Golden, and Assad [1990], and Gusfield [1988]). The U.S. Census Bureau produces a variety of tables from its census data. Suppose that it wishes to produce a $p \times q$ table $D = \{d_{ij}\}$ of nonnegative integers. Let $r(i)$ denote the sum of the matrix elements in the i th row and let $c(j)$ denote the sum of the matrix elements in the j th column. Assume that each sum $r(i)$ and $c(j)$ is strictly positive. The Census Bureau often wishes to disclose all the row and column sums along with some matrix elements (denoted by a set Y) and yet suppress the remaining elements to ensure the confidentiality of privileged information. Unless it exercises care, by disclosing the elements in Y , the Bureau might permit someone to deduce the exact value of one or more of the suppressed elements. It is possible to deduce a suppressed element d_{ij} if only one value of d_{ij} is consistent with the row and column sums and the disclosed elements in Y . We say that any such suppressed element is *unprotected*. Describe a polynomial-time algorithm for identifying all the unprotected elements of the matrix and their values.
- 6.6. **Minimax transportation problem** (Ahuja [1986]). Suppose that $G = (N, A)$ is an uncapacitated transportation problem (as defined in Section 1.2) and that we want to find an integer flow x that minimizes the objective function $\max\{c_{ij}x_{ij} : (i, j) \in A\}$ among all feasible integer flows.
- (a) Consider a relaxed version of the minimax transportation problem: Given a parameter λ , we want to know whether some feasible flow satisfies the condition $\max\{c_{ij}x_{ij} : (i, j) \in A\} \leq \lambda$. Show how to solve this problem as a maximum flow

- problem. [Hint: Use the condition $\max\{c_{ij}x_{ij} : (i, j) \in A\} \leq \lambda$ to formulate the problem as a feasible flow problem.]
- (b) Use the result in part (a) to develop a polynomial-time algorithm for solving the minimax transportation problem. What is the running time of your algorithm?
- 6.7. Consider a generalization of the feasible flow problem discussed in Application 6.1. Suppose that the flow bounds constraints are $l_{ij} \leq x_{ij} \leq u_{ij}$ instead of $0 \leq x_{ij} \leq u_{ij}$ for some nonnegative l_{ij} . How would you solve this generalization of the feasible flow problem as a single maximum flow problem?
- 6.8. Consider a generalization of the problem that we discussed in Application 6.2. Suppose that each club must nominate one of its members as a town representative so that the number of representatives belonging to the political party P_k is between l_k and u_k . Formulate this problem as a maximum flow problem with nonnegative lower bounds on arc flows.
- 6.9. In the example concerning the scheduling of uniform parallel machines (Application 6.4), we assumed that the same number of machines are available each day. How would you model a situation when the number of available machines varies from day to day? Illustrate your method on the example given in Application 6.4. Assume that three machines are available on days 1, 2, 4, and 5; two machines on days 3 and 6; and four machines on the rest of the days.
- 6.10. Can you solve the police patrol problem described in Exercise 1.9 using a maximum flow algorithm. If so, how?
- 6.11. Suppose that we wish to partition an undirected graph into two components with the minimum number of arcs between the components. How would you solve this problem?
- 6.12. Consider the network shown in Figure 6.22(a) together with the feasible flow x given in the figure.
- (a) Specify four s - t cuts in the network, each containing four forward arcs. List the capacity, residual capacity, and the flow across each cut.
- (b) Draw the residual network for the network given in Figure 6.22(a) and list four augmenting paths from node s to node t .

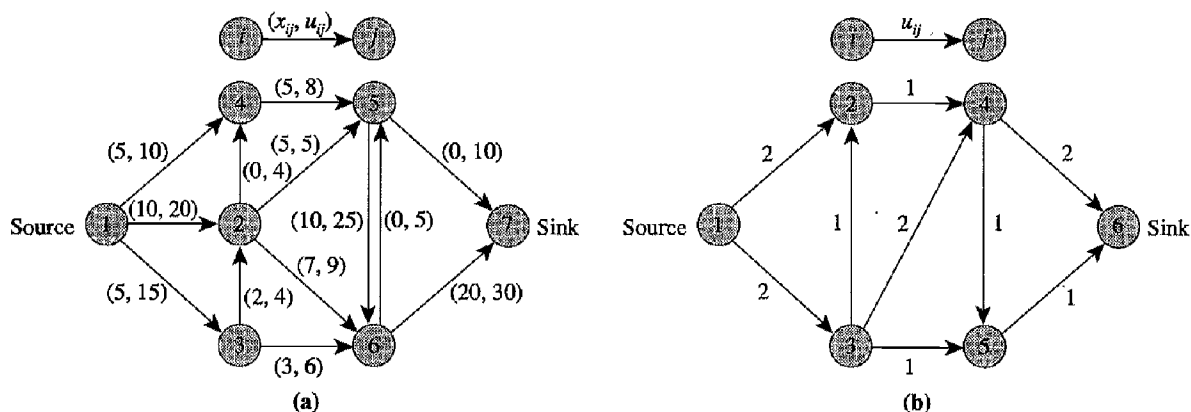


Figure 6.22 Examples for Exercises 6.12 and 6.13.

- 6.13. Solve the maximum flow problem shown in Figure 6.22(b) by the labeling algorithm, augmenting flow along the longest path in the residual network (i.e., the path containing maximum number of arcs). Specify the residual network before each augmentation. After every augmentation, decompose the flow into flows along directed paths from node s to node t . Finally, specify the minimum cut in the network obtained by the labeling algorithm.

- 6.14. Use the labeling algorithm to establish a maximum flow in the undirected network shown in Figure 6.23. Show the residual network at the end of each augmentation and specify the minimum cut that the algorithm obtains when it terminates.

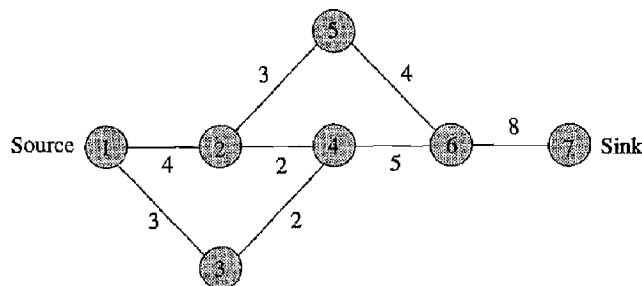


Figure 6.23 Example for Exercise 6.14.

- 6.15. Consider the network given in Figure 6.24; assume that each arc has capacity 1.
- Compute the maximum number of arc-disjoint paths from the source node to the sink node. (You might do so by inspection.)
 - Enumerate all s - t cuts in the network. For each s - t cut $[S, \bar{S}]$, list the node partition and the sets of forward and backward arcs.
 - Verify that the maximum number of arc-disjoint paths from node s to node t equals the minimum number of forward arcs in an s - t cut.

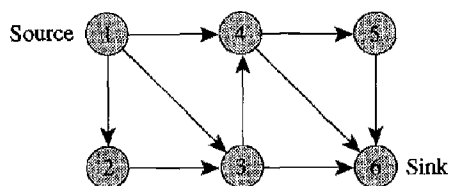


Figure 6.24 Example for Exercise 6.15.

- 6.16. Consider the matrix rounding problem given below (see Application 6.3). We want to round each element in the matrix, and also the row and column sums, to the nearest multiple of 2 so that the sum of the rounded elements in each row equals the rounded row sum and the sum of the rounded elements in each column equals the rounded column sum. Formulate this problem as a maximum flow problem and solve it.

	Row sum		
	7.5	6.3	15.4
	3.9	9.1	3.6
	15.0	5.5	21.5
Column sum	26.4	20.9	40.5

- 6.17. Formulate the following example of the scheduling problem on uniform parallel machines that we discussed in Application 6.4 as a maximum flow problem. Solve the problem by the labeling algorithm, assuming that two machines are available each day.

Job (j)	1	2	3	4
Processing time (p_j) (in days)	2.5	3.1	5.0	1.8
Release time (r_j)	1	5	0	2
Due date (d_j)	3	7	6	5

6.18. Minimum flow problem. The *minimum flow problem* is a close relative of the maximum flow problem with nonnegative lower bounds on arc flows. In the minimum flow problem, we wish to send the minimum amount of flow from the source to the sink, while satisfying given lower and upper bounds on arc flows.

(a) Show how to solve the minimum flow problem by using two applications of any maximum flow algorithm that applies to problems with zero lower bounds on arc flows. (*Hint*: First construct a feasible flow and then convert it into a minimum flow.)

(b) Prove the following *min-flow max-cut theorem*. Let the *floor* (or lower bound) on the cut capacity) of an s - t cut $[S, \bar{S}]$ be defined as $\sum_{(i,j) \in (S, \bar{S})} l_{ij} - \sum_{(i,j) \in (\bar{S}, S)} u_{ij}$. Show that the minimum value of all the flows from node s to node t equals the maximum floor of all s - t cuts.

6.19. Machine setup problem. A job shop needs to perform eight tasks on a particular day. Figure 6.25(a) shows the start and end times of each task. The workers must perform

Task	Start time	End time
1	1:00 P.M.	1:30 P.M.
2	6:00 P.M.	8:00 P.M.
3	10:00 P.M.	11:00 P.M.
4	4:00 P.M.	5:00 P.M.
5	4:00 P.M.	7:00 P.M.
6	12:00 noon	1:00 P.M.
7	2:00 P.M.	5:00 P.M.
8	11:00 P.M.	12:00 midnight

(a)

	1	2	3	4	5	6	7	8
1	—	60	10	25	30	20	15	40
2	10	—	40	55	40	5	30	35
3	65	30	—	0	45	30	20	5
4	0	50	35	—	20	15	10	20
5	20	24	40	50	—	15	5	23
6	10	8	9	35	12	—	30	30
7	15	30	6	18	15	30	—	10
8	20	35	15	12	75	13	25	—

(b)

Figure 6.25 Machine setup data: (a) task start and end times; (b) setup times in transforming between tasks.

these tasks according to this schedule so that exactly one worker performs each task. A worker cannot work on two jobs at the same time. Figure 6.25(b) shows the setup time (in minutes) required for a worker to go from one task to another. We wish to find the minimum number of workers to perform the tasks. Formulate this problem as a minimum flow problem (see Exercise 6.18).

- 6.20. Show how to transform a maximum flow problem having several source nodes and several sink nodes to one with only one source node and one sink node.
- 6.21. Show that if we add any number of incoming arcs, with any capacities, to the source node, the maximum flow value remains unchanged. Similarly, show that if we add any number of outgoing arcs, with any capacities, at the sink node, the maximum flow value remains unchanged.
- 6.22. Show that the maximum flow problem with integral data has a finite optimal solution if and only if the network contains no infinite capacity directed path from the source node to the sink node.
- 6.23. Suppose that a network has some infinite capacity arcs but no infinite capacity paths from the source to the sink. Let A^0 denote the set of arcs with finite capacities. Show that we can replace the capacity of each infinite capacity arc by a finite number $M \geq \sum_{(i,j) \in A^0} u_{ij}$ without affecting the maximum flow value.
- 6.24. Suppose that you want to solve a maximum flow problem containing parallel arcs, but the maximum flow code you own cannot handle parallel arcs. How would you use the code to solve your maximum flow problem?
- 6.25. **Networks with node capacities.** In some networks, in addition to arc capacities, each node i , other than the source and the sink, might have an upper bound, say $w(i)$, on the flow that can pass through it. For example, the nodes might be airports with limited runway capacity for takeoff and landings, or might be switches in a communication network with a limited number of ports. In these networks we are interested in determining the maximum flow satisfying both the arc and node capacities. Transform this problem to the standard maximum flow problem. From the perspective of worst-case complexity, is the maximum flow problem with upper bounds on nodes more difficult to solve than the standard maximum flow problem?
- 6.26. Suppose that a maximum flow network contains a node, other than the source node, with no incoming arc. Can we delete this node without affecting the maximum flow value? Similarly, can we delete a node, other than the sink node, with no outgoing arc?
- 6.27. Suppose that you are asked to solve a maximum flow problem in a directed network subject to the absolute value flow bound constraints $-u_{ij} \leq x_{ij} \leq u_{ij}$ imposed on some arcs (i, j) . How would you solve this problem?
- 6.28. Suppose that a maximum flow is available. Show how you would find a minimum cut in $O(m)$ additional time. Suppose, instead, that a minimum cut is available. Could you use this cut to obtain a maximum flow faster than applying a maximum flow algorithm?
- 6.29. **Painted network theorem.** Let G be a directed network with a distinguished arc (s, t) . Suppose that we paint each arc in the network as green, yellow, or red, with arc (s, t) painted yellow. Show that the painted network satisfies exactly one of the following two cases: (1) arc (s, t) is contained in a cycle of yellow and green arcs in which all yellow arcs have the same direction but green arcs can have arbitrary directions; (2) arc (s, t) is contained in a cut of yellow and red arcs in which all yellow arcs have the same direction but red arcs can have arbitrary directions.
- 6.30. Show that if $x_{ij} = u_{ij}$ for some arc (i, j) in every maximum flow, this arc must be a forward arc in some minimum cut.
- 6.31. An engineering department consisting of p faculty members, F_1, F_2, \dots, F_p , will offer p courses, C_1, C_2, \dots, C_p , in the coming semester and each faculty member will teach exactly one course. Each faculty member ranks two courses he (or she) would like to teach, ranking them according to his (or her) preference.
 - (a) We say that a course assignment is a *feasible* assignment if every faculty member

- teaches a course within his (or her) preference list. How would you determine whether the department can find a feasible assignment? (For a related problem see Exercise 12.46.)
- (b) A feasible assignment is said to be *k-feasible* if it assigns at most k faculty members to their second most preferred courses. For a given k , suggest an algorithm for determining a k -feasible assignment.
 - (c) We say that a feasible assignment is an *optimal assignment* if it maximizes the number of faculty members assigned to their most preferred course. Suggest an algorithm for determining an optimal assignment and analyze its complexity. [Hint: Use the algorithm in part (b) as a subroutine.]
- 6.32. **Airline scheduling problem.** An airline has p flight legs that it wishes to service by the fewest possible planes. To do so, it must determine the most efficient way to combine these legs into flight schedules. The starting time for flight i is a_i and the finishing time is b_i . The plane requires r_{ij} hours to return from the point of destination of flight i to the point of origin of flight j . Suggest a method for solving this problem.
- 6.33. A flow x is *even* if for every arc $(i, j) \in A$, x_{ij} is an even number; it is *odd* if for every $(i, j) \in A$, x_{ij} is an odd number. Either prove that each of the following claims are true or give a counterexample for them.
- (a) If all arc capacities are even, the network has an even maximum flow.
 - (b) If all arc capacities are odd, the network has an odd maximum flow.
- 6.34. Which of the following claims are true and which are false. Justify your answer either by giving a proof or by constructing a counterexample.
- (a) If x_{ij} is a maximum flow, either $x_{ij} = 0$ or $x_{ji} = 0$ for every arc $(i, j) \in A$.
 - (b) Any network always has a maximum flow x for which, for every arc $(i, j) \in A$, either $x_{ij} = 0$ or $x_{ji} = 0$.
 - (c) If all arcs in a network have different capacities, the network has a unique minimum cut.
 - (d) In a directed network, if we replace each directed arc by an undirected arc, the maximum flow value remains unchanged.
 - (e) If we multiply each arc capacity by a positive number λ , the minimum cut remains unchanged.
 - (f) If we add a positive number λ to each arc capacity, the minimum cut remains unchanged.
- 6.35. (a) Suppose that after solving a maximum flow problem you realize that you have underestimated the capacity of an arc (p, q) by k units. Show that the labeling algorithm can reoptimize the problem in $O(km)$ time.
- (b) Suppose that instead of underestimating the capacity of the arc (p, q) , you had overestimated its capacity by k units. Can you reoptimize the problem in $O(km)$ time?
- 6.36. (a) Construct a family of networks with the number of s - t cuts growing exponentially with n .
- (b) Construct a family of networks with the number of minimum cuts growing exponentially with n .
- 6.37. (a) Given a maximum flow in a network, describe an algorithm for determining the minimum cut $[S, \bar{S}]$ with the property that for every other minimum cut $[R, \bar{R}]$, $R \subseteq S$.
- (b) Describe an algorithm for determining the minimum cut $[S, \bar{S}]$ with the property that for every other minimum cut $[R, \bar{R}]$, $S \subseteq R$.
- (c) Describe an algorithm for determining whether the maximum flow problem has a unique minimum cut.
- 6.38. Let $[S, \bar{S}]$ and $[T, \bar{T}]$ be two s - t cuts in the directed network G . Show that the cut capacity function $u[\cdot, \cdot]$ is *submodular*, that is, $u[S, \bar{S}] + u[T, \bar{T}] \geq u[S \cup T, \overline{S \cup T}] + u[S \cap T, \overline{S \cap T}]$. (Hint: Prove this result by case analysis.)
- 6.39. Show that if $[S, \bar{S}]$ and $[T, \bar{T}]$ are both minimum cuts, so are $[S \cup T, \overline{S \cup T}]$ and $[S \cap T, \overline{S \cap T}]$.

- 6.40. Suppose that we know a noninteger maximum flow in a directed network with integer arc capacities. Suggest an algorithm for converting this flow into an integer maximum flow. What is the running time of your algorithm? (*Hint*: Send flows along cycles.)
- 6.41. **Optimal coverage of sporting events.** A group of reporters want to cover a set of sporting events in an olympiad. The sports events are held in several stadiums throughout a city. We know the starting time of each event, its duration, and the stadium where it is held. We are also given the travel times between different stadiums. We want to determine the least number of reporters required to cover the sporting events. How would you solve this problem?
- 6.42. In Section 6.7 we showed how to solve the maximum flow problem in directed networks with nonnegative lower bounds by solving two maximum flow problems with zero lower flow bounds. Try to generalize this approach for undirected networks in which the flow on any arc (i, j) is permitted in either direction, but whichever direction is chosen the amount of flow is at least l_{ij} . If you succeed in developing an algorithm, state the algorithm along with a proof that it correctly solves the problem; if you do not succeed in developing an algorithm state reasons why the generalization does not work.
- 6.43. **Feasibility of the feasible flow problem** (Gale [1957]). Show that the feasible flow problem, discussed in Application 6.1, has a feasible solution if and only if for every subset $S \subseteq N$, $b(S) - u[S, \bar{S}] \leq 0$. (*Hint*: Transform the feasible flow problem into a circulation problem with nonzero lower bounds and use the result of Theorem 6.11.)
- 6.44. Prove Theorems 6.7 and 6.8 for undirected networks.
- 6.45. Let N^+ and N^- be two nonempty disjoint node sets in G . Describe a method for determining the maximum number of arc-disjoint paths from N^+ to N^- (i.e., each path can start at any node in N^+ and can end at any node in N^-). What is the implication of the max-flow min-cut theorem in this case? (*Hint*: Generalize the statement of Theorem 6.7.)
- 6.46. Consider a 0–1 matrix H with n_1 rows and n_2 columns. We refer to a row or a column of the matrix H as a line. We say that a set of 1's in the matrix H is *independent* if no two of them appear in the same line. We also say that a set of lines in the matrix is a *cover* of H if they include (i.e., “cover”) all the 1's in the matrix. Show that the maximum number of independent 1's equals the minimum number of lines in a cover. (*Hint*: Use the max-flow min-cut theorem on an appropriately defined network.)
- 6.47. In a directed acyclic network G , certain arcs are colored blue, while others are colored red. Consider the problem of covering the blue arcs by directed paths, which can start and end at any node (these paths can contain arcs of any color). Show that the minimum number of directed paths needed to cover the blue arcs is equal to the maximum number of blue arcs that satisfy the property that no two of these arcs belong to the same path. Will this result be valid if G contains directed cycles? (*Hint*: Use the min-flow max-cut theorem stated in Exercise 6.18.)
- 6.48. **Pathological example for the labeling algorithm.** In the residual network $G(x)$ corresponding to a flow x , we define an *augmenting walk* as a directed walk from node s to node t that visits any arc at most once (it might visit nodes multiple times—in particular, an augmenting walk might visit nodes s and t multiple times.)
- Consider the network shown in Figure 6.26(a) with the arcs labeled a , b , c and d ; note that one arc capacity is irrational. Show that this network contains an infinite sequence of augmenting walks whose residual capacities sum to the maximum flow value. (*Hint*: Each augmenting walk of the sequence contains exactly two arcs from node s to node t with finite residual capacities.)
 - Now consider the network shown in Figure 6.26(b). Show that this network contains an infinite sequence of augmenting walks whose residual capacities sum to a value different than the maximum flow value.
 - Next consider the network shown in Figure 6.26(c); in addition to the arcs shown, the network contain an infinite capacity arc connecting each node pair in the set

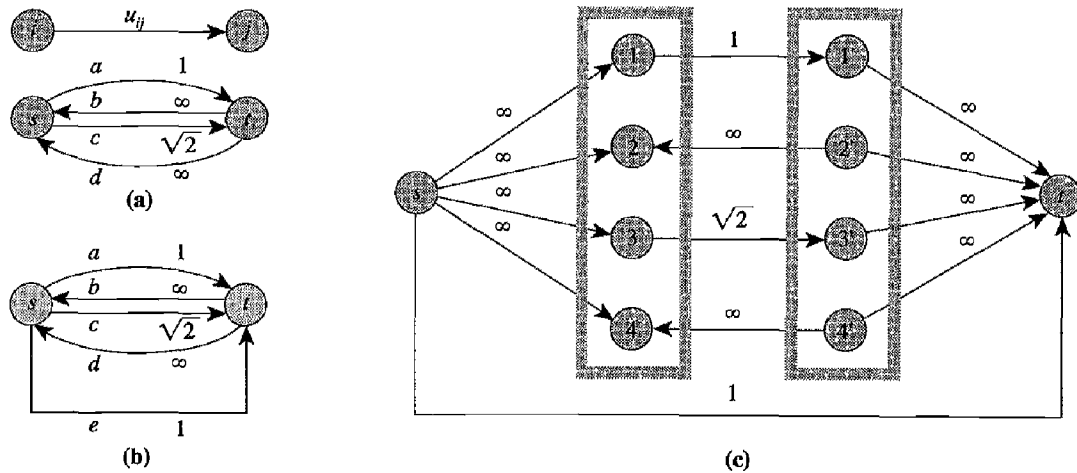


Figure 6.26 A subgraph of a pathological instance for labeling algorithm. The full graph contains an infinite capacity arc connecting each pair of nodes i and j as well as each pair of nodes i' and j' .

$\{1, 2, 3, 4\}$ and each node pair in the set $\{1', 2', 3', 4'\}$. Show that each augmenting walk in the solution of part (b) corresponds to an augmenting path in Figure 6.26(c). Conclude that the labeling algorithm, when applied to a maximum flow problem with irrational capacities, might perform an infinite sequence of augmentations and the terminal flow value might be different than the maximum flow value.

7

MAXIMUM FLOWS: POLYNOMIAL ALGORITHMS

Every day, in every way, I am getting better and better.
—Émile Coué

Chapter Outline

7.1	Introduction
7.2	Distance Labels
7.3	Capacity Scaling Algorithm
7.4	Shortest Augmenting Path Algorithm
7.5	Distance Labels and Layered Networks
7.6	Generic Preflow-Push Algorithm
7.7	FIFO Preflow-Push Algorithm
7.8	Highest Label Preflow-Push Algorithm
7.9	Excess Scaling Algorithm
7.10	Summary

7.1 INTRODUCTION

The generic augmenting path algorithm that we discussed in Chapter 6 is a powerful tool for solving maximum flow problems. Not only is it guaranteed to solve any maximum flow problem with integral capacity data, it also provides us with a constructive tool for establishing the fundamental max-flow min-cut theorem and therefore for developing many useful combinatorial applications of network flow theory.

As we noted in Chapter 6, however, the generic augmenting path algorithm has two significant computational limitations: (1) its worst-case computational complexity of $O(nmU)$ is quite unattractive for problems with large capacities; and (2) from a theoretical perspective, for problems with irrational capacity data, the algorithm might converge to a nonoptimal solution. These limitations suggest that the algorithm is not entirely satisfactory, in theory. Unfortunately, the algorithm is not very satisfactory in practice as well: On very large problems, it can require too much solution time.

Motivated by a desire to develop methods with improved worst-case complexity and empirical behavior, in this chapter we study several refinements of the generic augmenting path algorithm. We also introduce and study another class of algorithms, known as *preflow-push algorithms*, that have recently emerged as the most powerful techniques, both theoretically and computationally, for solving maximum flow problems.

Before describing these algorithms and analyzing them in detail, let us pause to reflect briefly on the theoretical limits of maximum flow algorithms and to introduce the general solution strategies employed by the refined augmenting path algorithms that we consider in this chapter. Flow decomposition theory shows that, in principle, we might be able to design augmenting path algorithms that are capable of finding a maximum flow in no more than m augmentations. For suppose that x is an optimal flow and x° is any initial flow (possibly the zero flow). By the flow decomposition property (see Section 3.5), we can obtain x from x° by a sequence of (1) at most m augmentations on augmenting paths from node s to node t , plus (2) flows around augmenting cycles. If we define x' as the flow vector obtained from x° by sending flows along only the augmenting paths, x' is also a maximum flow (because flows around augmenting cycles do not change the flow value into the sink node). This observation demonstrates a theoretical possibility of finding a maximum flow using at most m augmentations. Unfortunately, to apply this flow decomposition argument, we need to know a maximum flow. As a consequence, no algorithm developed in the literature achieves this theoretical bound of m augmentations. Nevertheless, it is possible to improve considerably on the $O(nU)$ bound on the number of augmentations required by the generic augmenting path algorithm.

How might we attempt to reduce the number of augmentations or even eliminate them altogether? In this chapter we consider three basic approaches:

1. Augmenting in “large” increments of flow
2. Using a combinatorial strategy that limits the type of augmenting paths we can use at each step
3. Relaxing the mass balance constraint at intermediate steps of the algorithm, and thus not requiring that each flow change must be an augmentation that starts at the source node and terminates at the sink node

Let us now consider each of these approaches. As we have seen in Chapter 6, the generic augmenting path algorithm could be slow because it might perform a large number of augmentations, each carrying a small amount of flow. This observation suggests one natural strategy for improving the augmenting path algorithm: Augment flow along a path with a *large* residual capacity so that the number of augmentations remains relatively *small*. The *maximum capacity augmenting path algorithm* uses this idea: It always augments flow along a path with the maximum residual capacity. In Section 7.3 we show that this algorithm performs $O(m \log U)$ augmentations. A variation of this algorithm that augments flows along a path with a *sufficiently large*, but not necessarily maximum residual capacity also performs $O(m \log U)$ augmentations and is easier to implement. We call this algorithm the *capacity scaling algorithm* and describe it in Section 7.3.

Another possible strategy for implementing and improving the augmenting path algorithm would be to develop an approach whose implementation is entirely independent of the arc capacity data and relies on a combinatorial argument for its convergence. One such approach would be somehow to restrict the choice of augmenting paths in some way. In one such approach we might always augment flow along a “shortest path” from the source to the sink, defining a shortest path as a directed path in the residual network consisting of the fewest number of arcs. If we

augment flow along a shortest path, the length of any shortest path either stays the same or increases. Moreover, within m augmentations, the length of the shortest path is guaranteed to increase. (We prove these assertions in Section 7.4.) Since no path contains more than $n - 1$ arcs, this result guarantees that the number of augmentations is at most $(n - 1)m$. We call this algorithm the *shortest augmenting path algorithm* and discuss it in Section 7.4.

The preflow-push algorithms use the third strategy we have identified: They seek out “shortest paths” as in the shortest augmenting path algorithm, but do not send flow along paths from the source to the sink. Instead, they send flows on individual arcs. This “localized” strategy, together with clever rules for implementing this strategy, permits these algorithms to obtain a speed-up not obtained by any augmenting path algorithm. We study these preflow-push algorithms in Sections 7.6 through 7.9.

The concept of distance labels is an important construct used to implement the shortest augmenting path algorithm and the preflow-push algorithms that we consider in this chapter. So before describing the improved algorithms, we begin by discussing this topic.

7.2 DISTANCE LABELS

A *distance function* $d: N \rightarrow Z^+ \cup \{0\}$ with respect to the residual capacities r_{ij} is a function from the set of nodes to the set of nonnegative integers. We say that a distance function is *valid* with respect to a flow x if it satisfies the following two conditions:

$$d(t) = 0; \tag{7.1}$$

$$d(i) \leq d(j) + 1 \quad \text{for every arc } (i, j) \text{ in the residual network } G(x). \tag{7.2}$$

We refer to $d(i)$ as the *distance label* of node i and conditions (7.1) and (7.2) as the *validity conditions*. The following properties show why the distance labels might be of use in designing network flow algorithms.

Property 7.1. *If the distance labels are valid, the distance label $d(i)$ is a lower bound on the length of the shortest (directed) path from node i to node t in the residual network.*

To establish the validity of this observation, let $i = i_1 - i_2 - \dots - i_k - i_{k+1} = t$ be any path of length k from node i to node t in the residual network. The validity conditions imply that

$$\begin{aligned} d(i_k) &\leq d(i_{k+1}) + 1 = d(t) + 1 = 1, \\ d(i_{k-1}) &\leq d(i_k) + 1 \leq 2, \\ d(i_{k-2}) &\leq d(i_{k-1}) + 1 \leq 3, \\ &\vdots \\ d(i) = d(i_1) &\leq d(i_2) + 1 \leq k. \end{aligned}$$

Property 7.2. If $d(s) \geq n$, the residual network contains no directed path from the source node to the sink node.

The correctness of this observation follows from the facts that $d(s)$ is a lower bound on the length of the shortest path from s to t in the residual network, and therefore no directed path can contain more than $(n - 1)$ arcs. Therefore, if $d(s) \geq n$, the residual network contains no directed path from node s to node t .

We now introduce some additional notation. We say that the distance labels are *exact* if for each node i , $d(i)$ equals the length of the shortest path from node i to node t in the residual network. For example, in Figure 7.1, if node 1 is the source node and node 4 is the sink node, then $d = (0, 0, 0, 0)$ is a valid vector of distance labels, and $d = (3, 1, 2, 0)$ is a vector of exact distance labels. We can determine exact distance labels for all nodes in $O(m)$ time by performing a backward breadth-first search of the network starting at the sink node (see Section 3.4).

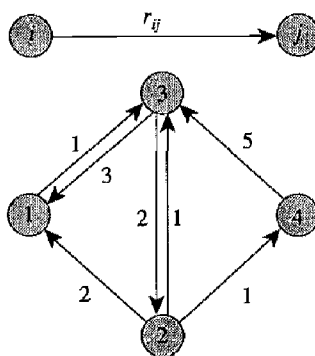


Figure 7.1 Residual network.

Admissible Arcs and Admissible Paths

We say that an arc (i, j) in the residual network is *admissible* if it satisfies the condition that $d(i) = d(j) + 1$; we refer to all other arcs as *inadmissible*. We also refer to a path from node s to node t consisting entirely of admissible arcs as an *admissible path*. Later, we use the following property of admissible paths.

Property 7.3. An admissible path is a shortest augmenting path from the source to the sink.

Since every arc (i, j) in an admissible path P is admissible, the residual capacity of this arc and the distance labels of its end nodes satisfy the conditions (1) $r_{ij} > 0$, and (2) $d(i) = d(j) + 1$. Condition (1) implies that P is an augmenting path and condition (2) implies that if P contains k arcs, then $d(s) = k$. Since $d(s)$ is a lower bound on the length of any path from the source to the sink in the residual network (from Property 7.1), the path P must be a shortest augmenting path.

7.3 CAPACITY SCALING ALGORITHM

We begin by describing the maximum capacity augmenting path algorithm and noting its computational complexity. This algorithm always augments flow along a path with the maximum residual capacity. Let x be any flow and let v be its flow value.

As before, let v^* be the maximum flow value. The flow decomposition property (i.e., Theorem 3.5), as applied to the residual network $G(x)$, implies that we can find m or fewer directed paths from the source to the sink whose residual capacities sum to $(v^* - v)$. Thus the maximum capacity augmenting path has residual capacity at least $(v^* - v)/m$. Now consider a sequence of $2m$ consecutive maximum capacity augmentations starting with the flow x . If each of these augmentations augments at least $(v^* - v)/2m$ units of flow, then within $2m$ or fewer iterations we will establish a maximum flow. Note, however, that if one of these $2m$ consecutive augmentations carries less than $(v^* - v)/2m$ units of flow, then from the initial flow vector x , we have reduced the residual capacity of the maximum capacity augmenting path by a factor of at least 2. This argument shows that within $2m$ consecutive iterations, the algorithm either establishes a maximum flow or reduces the residual capacity of the maximum capacity augmenting path by a factor of at least 2. Since the residual capacity of any augmenting path is at most $2U$ and is at least 1, after $O(m \log U)$ iterations, the flow must be maximum. (Note that we are essentially repeating the argument used to establish the geometric improvement approach discussed in Section 3.3.)

As we have seen, the maximum capacity augmentation algorithm reduces the number of augmentations in the generic labeling algorithm from $O(nU)$ to $O(m \log U)$. However, the algorithm performs more computations per iteration, since it needs to identify an augmenting path with the maximum residual capacity, not just any augmenting path. We now suggest a variation of the maximum capacity augmentation algorithm that does not perform more computations per iteration and yet establishes a maximum flow within $O(m \log U)$. Since this algorithm scales the arc capacities implicitly, we refer to it as the *capacity scaling algorithm*.

The essential idea underlying the capacity scaling algorithm is conceptually quite simple: We augment flow along a path with a *sufficiently large* residual capacity, instead of a path with the maximum augmenting capacity because we can obtain a path with a sufficiently large residual capacity fairly easily—in $O(m)$ time. To define the capacity scaling algorithm, let us introduce a parameter Δ and, with respect to a given flow x , define the Δ -residual network as a network containing arcs whose residual capacity is at least Δ . Let $G(x, \Delta)$ denote the Δ -residual network. Note that $G(x, 1) = G(x)$ and $G(x, \Delta)$ is a subgraph of $G(x)$. Figure 7.2 illustrates this definition. Figure 7.2(a) gives the residual network $G(x)$ and Figure 7.2(b) gives the Δ -residual network $G(x, \Delta)$ for $\Delta = 8$. Figure 7.3 specifies the capacity scaling algorithm.

Let us refer to a phase of the algorithm during which Δ remains constant as a *scaling phase* and a scaling phase with a specific value of Δ as a Δ -scaling phase. Observe that in a Δ -scaling phase, each augmentation carries at least Δ units of flow. The algorithm starts with $\Delta = 2^{\lceil \log U \rceil}$ and halves its value in every scaling phase until $\Delta = 1$. Consequently, the algorithm performs $1 + \lceil \log U \rceil = O(\log U)$ scaling phases. In the last scaling phase, $\Delta = 1$, so $G(x, \Delta) = G(x)$. This result shows that the algorithm terminates with a maximum flow.

The efficiency of the algorithm depends on the fact that it performs at most $2m$ augmentations per scaling phase. To establish this result, consider the flow at the end of the Δ -scaling phase. Let x' be this flow and let v' denote its flow value. Furthermore, let S be the set of nodes reachable from node s in $G(x', \Delta)$. Since

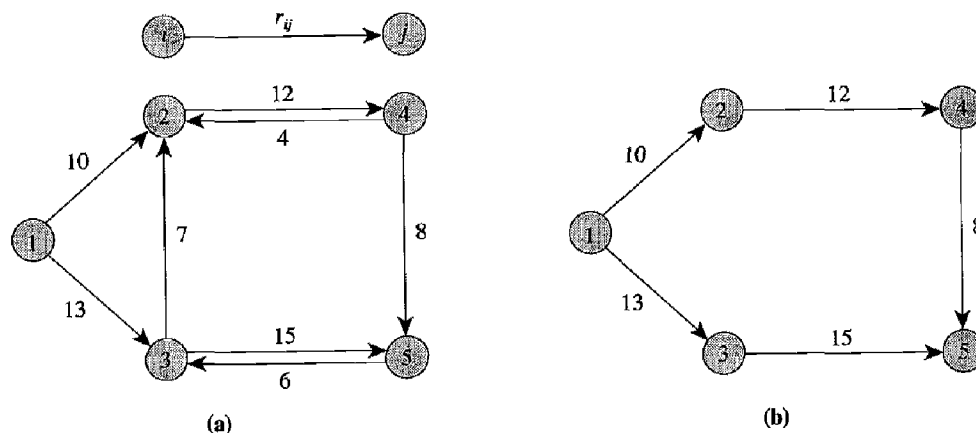


Figure 7.2 Illustrating the Δ -residual network: (a) residual network $G(x)$; (b) Δ -residual network $G(x, \Delta)$ for $\Delta = 8$.

```

algorithm capacity scaling;
begin
   $x := 0$ ;
   $\Delta := 2^{\lceil \log U \rceil}$ ;
  while  $\Delta \geq 1$  do
    begin
      while  $G(x, \Delta)$  contains a path from node  $s$  to node  $t$  do
        begin
          identify a path  $P$  in  $G(x, \Delta)$ ;
           $\delta := \min\{r_{ij} : (i, j) \in P\}$ ;
          augment  $\delta$  units of flow along  $P$  and update  $G(x, \Delta)$ ;
        end;
         $\Delta := \Delta/2$ ;
      end;
    end;
end;

```

Figure 7.3 Capacity scaling algorithm.

$G(x', \Delta)$ contains no augmenting path from the source to the sink, $t \in S$. Therefore, $[S, \bar{S}]$ forms an s - t cut. The definition of S implies that the residual capacity of every arc in $[S, \bar{S}]$ is strictly less than Δ , so the residual capacity of the cut $[S, \bar{S}]$ is at most $m\Delta$. Consequently, $v^* - v' \leq m\Delta$ (from Property 6.2). In the next scaling phase, each augmentation carries at least $\Delta/2$ units of flow, so this scaling phase can perform at most $2m$ such augmentations. The labeling algorithm described in Section 6.5 requires $O(m)$ time to identify an augmenting path, and updating the Δ -residual network also requires $O(m)$ time. These arguments establish the following result.

Theorem 7.4. *The capacity scaling algorithm solves the maximum flow problem within $O(m \log U)$ augmentations and runs in $O(m^2 \log U)$ time.* ♦

It is possible to reduce the complexity of the capacity scaling algorithm even further—to $O(nm \log U)$ —using ideas of the shortest augmenting path algorithm, described in the next section.

7.4 SHORTEST AUGMENTING PATH ALGORITHM

The shortest augmenting path algorithm always augments flow along a shortest path from the source to the sink in the residual network. A natural approach for implementing this approach would be to look for shortest paths by performing a breadth first search in the residual network. If the labeling algorithm maintains the set L of labeled nodes as a queue, then by examining the labeled nodes in a first-in, first-out order, it would obtain a shortest path in the residual network (see Exercise 3.30). Each of these iterations would require $O(m)$ steps in the worst case, and (by our subsequent observations) the resulting computation time would be $O(nm^2)$. Unfortunately, this computation time is excessive. We can improve it by exploiting the fact that the minimum distance from any node i to the sink node t is monotonically nondecreasing over all augmentations. By fully exploiting this property, we can reduce the average time per augmentation to $O(n)$.

The shortest augmenting path algorithm proceeds by augmenting flows along admissible paths. It constructs an admissible path incrementally by adding one arc at a time. The algorithm maintains a *partial admissible path* (i.e., a path from s to some node i consisting solely of admissible arcs) and iteratively performs *advance* or *retreat* operations from the last node (i.e., the tip) of the partial admissible path, which we refer to as the *current node*. If the current node i has (i.e., is incident to) an admissible arc (i, j) , we perform an advance operation and add arc (i, j) to the partial admissible path; otherwise, we perform a retreat operation and backtrack one arc. We repeat these operations until the partial admissible path reaches the sink node at which time we perform an augmentation. We repeat this process until the flow is maximum. Before presenting a formal description of the algorithm, we illustrate it on the numerical example given in Figure 7.4(a).

We first compute the initial distance labels by performing the backward breadth-first search of the residual network starting at the sink node. The numbers next to the nodes in Figure 7.4(a) specify these values of the distance labels. In this example we adopt the convention of selecting the arc (i, j) with the smallest value of j whenever node i has several admissible arcs. We start at the source node with a null partial admissible path. The source node has several admissible arcs, so we perform an advance operation. This operation adds the arc $(1, 2)$ to the partial admissible path. We store this path using predecessor indices, so we set $\text{pred}(2) = 1$. Now node 2 is the current node and the algorithm performs an advance operation at node 2. In doing so, it adds arc $(2, 7)$ to the partial admissible path, which now becomes $1-2-7$. We also set $\text{pred}(7) = 2$. In the next iteration, the algorithm adds arc $(7, 12)$ to the partial admissible path obtaining $1-2-7-12$, which is an admissible path to the sink node. We perform an augmentation of value $\min\{r_{12}, r_{27}, r_{7,12}\} = \min\{2, 1, 2\} = 1$, and thus saturate the arc $(2, 7)$. Figure 7.4(b) specifies the residual network at this stage.

We again start at the source node with a null partial admissible path. The algorithm adds the arc $(1, 2)$ and node 2 becomes the new current node. Now we find that node 2 has no admissible arc. To create new admissible arcs, we must increase the distance label of node 2. We thus increase $d(2)$ to the value $\min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\} = \min\{d(1) + 1\} = 4$. We refer to this operation as a *relabel* operation. We will later show that a relabel operation preserves the validity

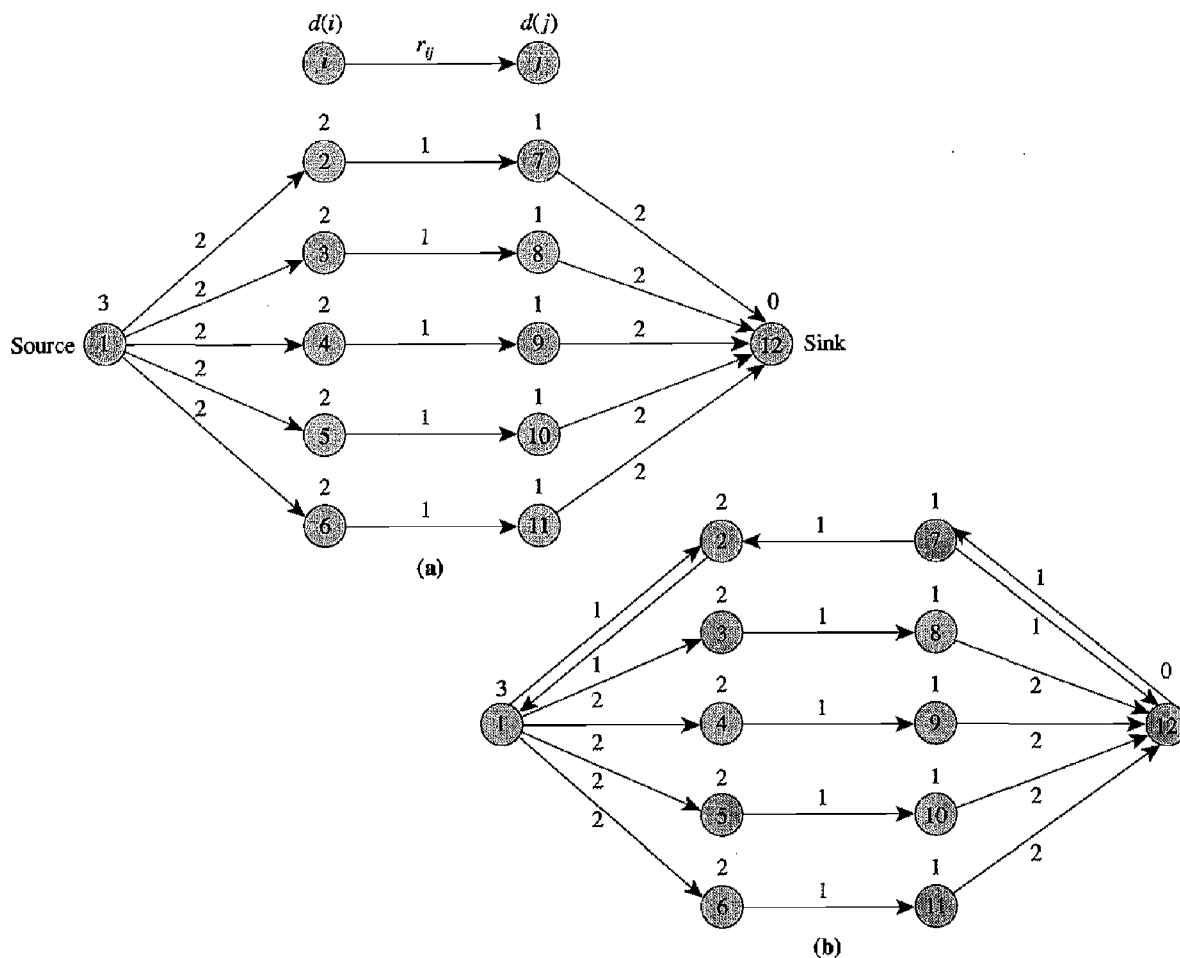


Figure 7.4 Illustrating the shortest augmenting path algorithm.

conditions imposed upon the distance labels. Observe that the increase in $d(2)$ causes arc $(1, 2)$ to become inadmissible. Thus we delete arc $(1, 2)$ from the partial admissible path which again becomes a null path. In the subsequent operations, the algorithm identifies the admissible paths $1-3-8-12$, $1-4-9-12$, $1-5-10-12$, and $1-6-11-12$ and augments unit flows on these paths. We encourage the reader to carry out the details of these operations. Figures 7.5 and 7.6 specify the details of the algorithm.

Correctness of the Algorithm

In our analysis of the shortest augmenting path algorithm we first show that it correctly solves the maximum flow problem.

Lemma 7.5. *The shortest augmenting path algorithm maintains valid distance labels at each step. Moreover, each relabel (or, retreat) operation strictly increases the distance label of a node.*

```

algorithm shortest augmenting path;
begin
   $x := 0$ ;
  obtain the exact distance labels  $d(i)$ ;
   $i := s$ ;
  while  $d(s) < n$  do
    begin
      if  $i$  has an admissible arc then
        begin
          advance( $i$ );
          if  $i = t$  then augment and set  $i = s$ 
        end
      else retreat( $i$ )
    end;
  end;

```

Figure 7.5 Shortest augmenting path algorithm.

```

procedure advance( $i$ );
begin
  let  $(i, j)$  be an admissible arc in  $A(i)$ ;
   $\text{pred}(j) := i$  and  $i := j$ ;
end;

```

(a)

```

procedure retreat( $i$ );
begin
   $d(i) := \min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\}$ ;
  if  $i \neq s$  then  $i := \text{pred}(i)$ ;
end;

```

(b)

```

procedure augment;
begin
  using the predecessor indices identify an augmenting
  path  $P$  from the source to the sink;
   $\delta := \min\{r_{ij} : (i, j) \in P\}$ ;
  augment  $\delta$  units of flow along path  $P$ ;
end;

```

(c)

Figure 7.6 Procedures of the shortest augmenting path algorithm.

Proof. We show that the algorithm maintains valid distance labels at every step by performing induction on the number of augment and relabel operations. (The advance operation does not affect the admissibility of any arc because it does not change any residual capacity or distance label.) Initially, the algorithm constructs valid distance labels. Assume, inductively, that the distance labels are valid prior to an operation (i.e., they satisfy the validity conditions). We need to check whether these conditions remain valid (a) after an augment operation, and (b) after a relabel operation.

(a) Although a flow augmentation on arc (i, j) might remove this arc from the residual network, this modification to the residual network does not affect the validity of the distance labels for this arc. An augmentation on arc (i, j) might, however, create an additional arc (j, i) with $r_{ji} > 0$ and therefore also create an additional inequality $d(j) \leq d(i) + 1$ that the distance labels must satisfy. The distance labels satisfy this validity condition, though, since $d(i) = d(j) + 1$ by the admissibility property of the augmenting path.

(b) The relabel operation modifies $d(i)$; therefore, we must show that each incoming and outgoing arc at node i satisfies the validity conditions with respect to the new distance labels, say $d'(i)$. The algorithm performs a relabel operation at node i when it has no admissible arc; that is, no arc $(i, j) \in A(i)$ satisfies the conditions $d(i) = d(j) + 1$ and $r_{ij} > 0$. This observation, in light of the validity condition $d(i) \leq d(j) + 1$, implies that $d(i) < d(j) + 1$ for all arcs $(i, j) \in A$ with a positive residual capacity. Therefore, $d(i) < \min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\} = d'(i)$, which is the new distance label after the relabel operation. We have thus shown that relabeling preserves the validity condition for all arcs emanating from node i , and that each relabel operation strictly increases the value of $d(i)$. Finally, note that every incoming arc (k, i) satisfies the inequality $d(k) \leq d(i) + 1$ (by the induction hypothesis). Since $d(i) < d'(i)$, the relabel operation again preserves validity condition for arc (k, i) . ♦

The shortest augmenting path algorithm terminates when $d(s) \geq n$, indicating that the network contains no augmenting path from the source to the sink (from Property 7.2). Consequently, the flow obtained at the end of the algorithm is a maximum flow. We have thus proved the following theorem.

Theorem 7.6. *The shortest augmenting path algorithm correctly computes a maximum flow.* ♦

Complexity of the Algorithm

We now show that the shortest augmenting path algorithm runs in $O(n^2m)$ time. We first describe a data structure used to select an admissible arc emanating from a given node. We call this data structure the *current-arc data structure*. Recall that we used this data structure in Section 3.4 in our discussion of search algorithms. We also use this data structure in almost all the maximum flow algorithms that we describe in subsequent sections. Therefore, we review this data structure before proceeding.

Recall that we maintain the arc list $A(i)$ which contains all the arcs emanating from node i . We can arrange the arcs in these lists arbitrarily, but the order, once decided, remains unchanged throughout the algorithm. Each node i has a *current arc*, which is an arc in $A(i)$ and is the next candidate for admissibility testing. Initially, the current arc of node i is the first arc in $A(i)$. Whenever the algorithm attempts to find an admissible arc emanating from node i , it tests whether the node's current arc is admissible. If not, it designates the next arc in the arc list as the current arc. The algorithm repeats this process until either it finds an admissible arc or reaches the end of the arc list.

Consider, for example, the arc list of node 1 in Figure 7.7. In this instance, $A(1) = \{(1, 2), (1, 3), (1, 4), (1, 5), (1, 6)\}$. Initially, the current arc of node 1 is arc (1, 2). Suppose that the algorithm attempts to find an admissible arc emanating from node 1. It checks whether the node's current arc, arc (1, 2), is admissible. Since it is not, the algorithm designates arc (1, 3) as the current arc of node 1. The arc (1, 3) is also inadmissible, so the current arc becomes arc (1, 4), which is admissible. From this point on, arc (1, 4) remains the current arc of node 1 until it becomes inadmissible because the algorithm has increased the value of $d(4)$ or decreased the value of the residual capacity of arc (1, 4) to zero.

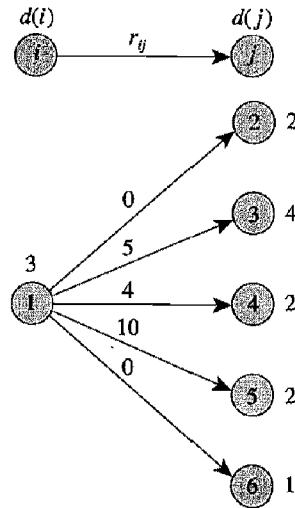


Figure 7.7 Selecting admissible arcs emanating from a node.

Let us consider the situation when the algorithm reaches the end of the arc list without finding any admissible arc. Can we say that $A(i)$ has no admissible arc? We can, because it is possible to show that if an arc (i, j) is inadmissible in previous iterations, it remains inadmissible until $d(i)$ increases (see Exercise 7.13). So if we reach the end of the arc list, we perform a relabel operation and again set the current arc of node i to be the first arc in $A(i)$. The relabel operation also examines each arc in $A(i)$ once to compute the new distance label, which is same as the time it spends in identifying admissible arcs at node i in one scan of the arc list. We have thus established the following result.

Property 7.7. *If the algorithm relabels any node at most k times, the total time spent in finding admissible arcs and relabeling the nodes is $O(k \sum_{i \in N} |A(i)|) = O(km)$.*

We shall be using this result several times in this and the following chapters. We also use the following result in several places.

Lemma 7.8. *If the algorithm relabels any node at most k times, the algorithm saturates arcs (i.e., reduces their residual capacity to zero) at most $km/2$ times.*

Proof. We show that between two consecutive saturations of an arc (i, j) , both $d(i)$ and $d(j)$ must increase by at least 2 units. Since, by our hypothesis, the algorithm increases each distance label at most k times, this result would imply that the algorithm could saturate any arc at most $k/2$ times. Therefore, the total number of arc saturations would be $km/2$, which is the assertion of the lemma.

Suppose that an augmentation saturates an arc (i, j) . Since the arc (i, j) is admissible,

$$d(i) = d(j) + 1. \quad (7.3)$$

Before the algorithm saturates this arc again, it must send flow back from node j to node i . At this time, the distance labels $d'(i)$ and $d'(j)$ satisfy the equality

$$d'(j) = d'(i) + 1. \quad (7.4)$$

In the next saturation of arc (i, j) , we must have

$$d''(i) = d''(j) + 1. \quad (7.5)$$

Using (7.3) and (7.4) in (7.5), we see that

$$d''(i) = d''(j) + 1 \geq d'(j) + 1 = d'(i) + 2 \geq d(i) + 2.$$

The inequalities in this expression follow from Lemma 7.5. Similarly, it is possible to show that $d''(j) \geq d(j) + 2$. As a result, between two consecutive saturations of the arc (i, j) , both $d(i)$ and $d(j)$ increase by at least 2 units, which is the conclusion of the lemma. ♦

Lemma 7.9.

- (a) *In the shortest augmenting path algorithm each distance label increases at most n times. Consequently, the total number of relabel operations is at most n^2 .*
- (b) *The number of augment operations is at most $nm/2$.*

Proof. Each relabel operation at node i increases the value of $d(i)$ by at least 1 unit. After the algorithm has relabeled node i at most n times, $d(i) \geq n$. From this point on, the algorithm never again selects node i during an advance operation since for every node k in the partial admissible path, $d(k) < d(s) < n$. Thus the algorithm relabels a node at most n times and the total number of relabel operations is bounded by n^2 . In view of Lemma 7.8, the preceding result implies that the algorithm saturates at most $nm/2$ arcs. Since each augmentation saturates at least one arc, we immediately obtain a bound of $nm/2$ on the number of augmentations. ♦

Theorem 7.10. *The shortest augmenting path algorithm runs in $O(n^2m)$ time.*

Proof. Using Lemmas 7.9 and 7.7 we find that the total effort spent in finding admissible arcs and in relabeling the nodes is $O(nm)$. Lemma 7.9 implies that the total number of augmentations is $O(nm)$. Since each augmentation requires $O(n)$ time, the total effort for the augmentation operations is $O(n^2m)$. Each retreat operation relabels a node, so the total number of retreat operations is $O(n^2)$. Each advance operation adds one arc to the partial admissible path, and each retreat operation deletes one arc from it. Since each partial admissible path has length at most n , the algorithm requires at most $O(n^2 + n^2m)$ advance operations. The first

term comes from the number of retreat (relabel) operations, and the second term from the number of augmentations. The combination of these bounds establishes the theorem. ♦

A Practical Improvement

The shortest augmenting path algorithm terminates when $d(s) \geq n$. This termination criteria is satisfactory for the worst-case analysis but might not be efficient in practice. Empirical investigations have revealed that the algorithm spends too much time relabeling nodes and that a major portion of this effort is performed after the algorithm has established a maximum flow. This happens because the algorithm does not know that it has found a maximum flow. We next suggest a technique that is capable of detecting the presence of a minimum cut and so the existence of a maximum flow much before the label of node s satisfies the condition $d(s) \geq n$. Incorporating this technique in the shortest augmenting path algorithm improves its performance substantially in practice.

We illustrate this technique by applying it to the numerical example we used earlier to illustrate the shortest augmenting path algorithm. Figure 7.8 gives the residual network immediately after the last augmentation. Although the flow is now a maximum flow, since the source is not connected to the sink in the residual network, the termination criteria of $d(1) \geq 12$ is far from being satisfied. The reader can verify that after the last augmentation, the algorithm would increase the distance labels of nodes 6, 1, 2, 3, 4, 5, in the given order, each time by 2 units. Eventually, $d(1) \geq 12$ and the algorithm terminates. Observe that the node set S of the minimum cut $[S, \bar{S}]$ equals $\{6, 1, 2, 3, 4, 5\}$, and the algorithm increases the distance labels of all the nodes in S without performing any augmentation. The technique we describe essentially detects a situation like this one.

To implement this approach, we maintain an n -dimensional additional array,

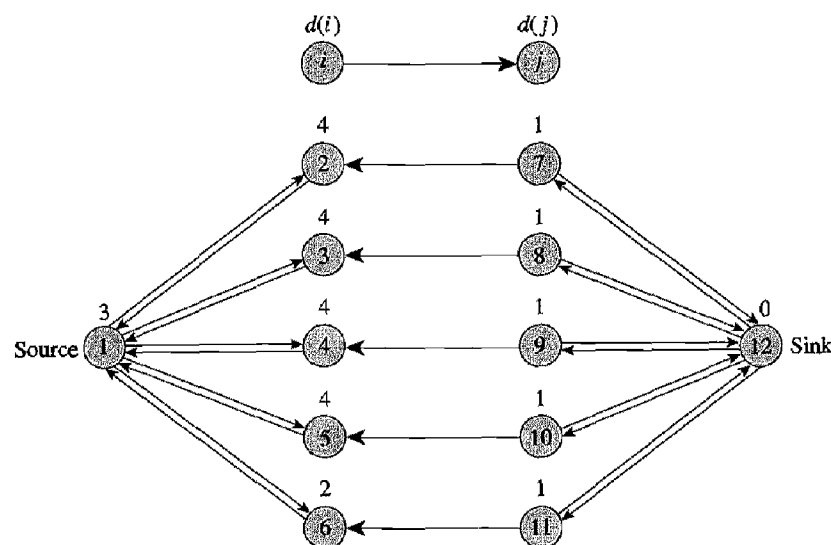


Figure 7.8 Bad example for the shortest augmenting path algorithm.

numb, whose indices vary from 0 to $(n - 1)$. The value $\text{numb}(k)$ is the number of nodes whose distance label equals k . The algorithm initializes this array while computing the initial distance labels using a breadth first search. At this point, the positive entries in the array *numb* are consecutive [i.e., the entries $\text{numb}(0)$, $\text{numb}(1)$, \dots , $\text{numb}(l)$ will be positive up to some index l and the remaining entries will all be zero]. For example, the *numb* array for the distance labels shown in Figure 7.8 is $\text{numb}(0) = 1$, $\text{numb}(1) = 5$, $\text{numb}(2) = 1$, $\text{numb}(3) = 1$, $\text{numb}(4) = 4$ and the remaining entries are zero. Subsequently, whenever the algorithm increases the distance label of a node from k_1 to k_2 , it subtracts 1 from $\text{numb}(k_1)$, adds 1 to $\text{numb}(k_2)$ and checks whether $\text{numb}(k_1) = 0$. If $\text{numb}(k_1)$ does equal zero, the algorithm terminates. As seen earlier, the shortest augmenting path algorithm augments unit flow along the paths 1–2–7–12, 1–3–8–12, 1–4–9–12, 1–5–10–12, and 1–6–11–12. At the end of these augmentations, we obtain the residual network shown in Figure 7.8. When we continue the shortest augmenting path algorithm from this point, it constructs the partial admissible path 1–6. Next it relabels node 6 and its distance label increases from 2 to 4. The algorithm finds that $\text{numb}(2) = 0$ and it terminates.

To see why this termination criterion works, let $S = \{i \in N : d(i) > k_1\}$ and $\bar{S} = \{i \in N : d(i) < k_1\}$. It is easy to verify that $s \in S$ and $t \in \bar{S}$. Now consider the s – t cut $[S, \bar{S}]$. The definitions of the sets S and \bar{S} imply that $d(i) > d(j) + 1$ for all $(i, j) \in [S, \bar{S}]$. The validity condition (7.2) implies that $r_{ij} = 0$ for each arc $(i, j) \in [S, \bar{S}]$. Therefore, $[S, \bar{S}]$ is a minimum cut and the current flow is a maximum flow.

Application to Capacity Scaling Algorithm

In the preceding section we described an $O(m^2 \log U)$ time capacity scaling algorithm for the maximum flow problem. We can improve the running time of this algorithm to $O(nm \log U)$ by using the shortest augmenting path as a subroutine in the capacity scaling algorithm. Recall that the capacity scaling algorithm performs a number of Δ -scaling phases and in the Δ -scaling phase sends the maximum possible flow in the Δ -residual network $G(x, \Delta)$, using the labeling algorithm as a subroutine. In the improved implementation, we use the shortest augmenting path algorithm to send the maximum possible flow from node s to node t . We accomplish this by defining the distance labels with respect to the network $G(x, \Delta)$ and augmenting flow along the shortest augmenting path in $G(x, \Delta)$. Recall from the preceding section that a scaling phase contains $O(m)$ augmentations. The complexity analysis of the shortest augmenting path algorithm implies that if the algorithm is guaranteed to perform $O(m)$ augmentations, it would run in $O(nm)$ time because the time for augmentations reduces from $O(n^2m)$ to $O(nm)$ and all other operations, as before, require $O(nm)$ time. These observations immediately yield a bound of $O(nm \log U)$ on the running time of the capacity scaling algorithm.

Further Worst-Case Improvements

The idea of augmenting flows along shortest paths is intuitively appealing and easy to implement in practice. The resulting algorithms identify at most $O(nm)$ augmenting paths and this bound is tight [i.e., on particular examples these algorithms perform $\Omega(nm)$ augmentations]. The only way to improve the running time of the shortest

augmenting path algorithm is to perform fewer computations per augmentation. The use of a sophisticated data structure, called *dynamic trees*, reduces the average time for each augmentation from $O(n)$ to $O(\log n)$. This implementation of the shortest augmenting path algorithm runs in $O(nm \log n)$ time, and obtaining further improvements appears quite difficult except in very dense networks. We describe the dynamic tree implementation of the shortest augmenting path algorithm in Section 8.5.

7.5 DISTANCE LABELS AND LAYERED NETWORKS

Like the shortest augmenting path algorithm, several other maximum flow algorithms send flow along shortest paths from the source to the sink. Dinic's algorithm is a popular algorithm in this class. This algorithm constructs shortest path networks, called *layered networks*, and establishes *blocking flows* (to be defined later) in these networks. In this section we point out the relationship between layered networks and distance labels. By developing a modification of the shortest augmenting path algorithm that reduces to Dinic's algorithm, we show how to use distance labels to simulate layered networks.

With respect to a given flow x , we define the *layered network* V as follows. We determine the exact distance labels d in $G(x)$. The layered network consists of those arcs (i, j) in $G(x)$ satisfying the condition $d(i) = d(j) + 1$. For example, consider the residual network $G(x)$ given in Figure 7.9(a). The number beside each node represents its exact distance label. Figure 7.9(b) shows the layered network of $G(x)$. Observe that by definition every path from the source to the sink in the layered network V is a shortest path in $G(x)$. Observe further that some arc in V might not be contained in any path from the source to the sink. For example, in Figure 7.9(b), arcs $(5, 7)$ and $(6, 7)$ do not lie on any path in V from the source to the sink. Since these arcs do not participate in any flow augmentation, we typically delete them from the layered network; doing so gives us Figure 7.9(c). In the resulting layered network, the nodes are partitioned into layers of nodes $V_0, V_1, V_2, \dots, V_i$; layer k contains the nodes whose distance labels equal k . Furthermore, for every

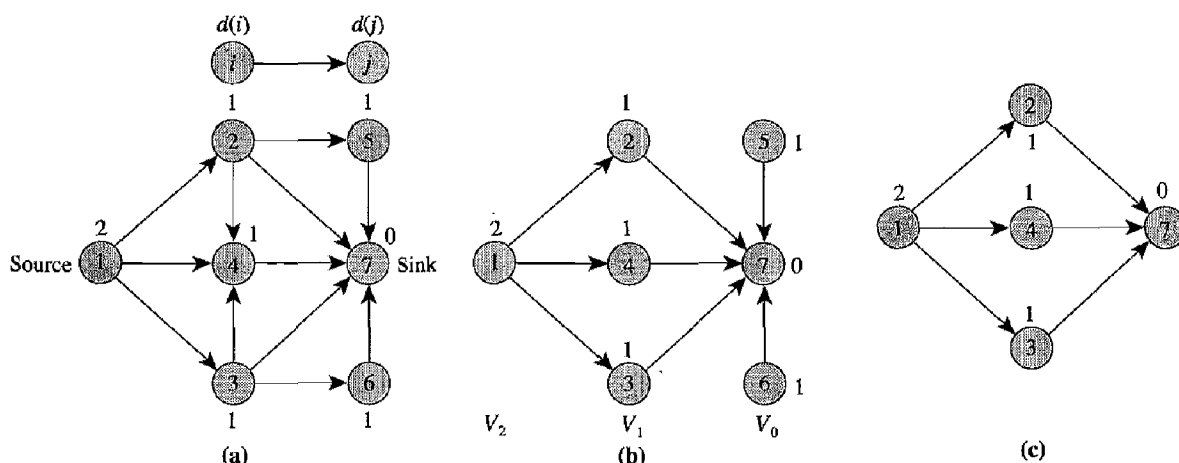


Figure 7.9 Forming layered networks: (a) residual network; (b) corresponding layered network; (c) layered network after deleting redundant arcs.

arc (i, j) in the layered network, $i \in V_k$ and $j \in V_{k-1}$ for some k . Let the source node have the distance label l .

Dinic's algorithm proceeds by augmenting flows along directed paths from the source to the sink in the layered network. The augmentation of flow along an arc (i, j) reduces the residual capacity of arc (i, j) and increases the residual capacity of the reversal arc (j, i) ; however, each arc of the layered network is admissible, and therefore Dinic's algorithm does not add reversal arcs to the layered network. Consequently, the length of every augmenting path is $d(s)$ and in an augmenting path every arc (i, j) has $i \in V_k$ and $j \in V_{k-1}$ for some k . The latter fact allows us to determine an augmenting path in the layered network, on average, in $O(n)$ time. (The argument used to establish the $O(n)$ time bound is the same as that used in our analysis of the shortest augmenting path algorithm.) Each augmentation saturates at least one arc in the layered network, and after at most m augmentations the layered network contains no augmenting path. We call the flow at this stage a *blocking flow*. We have shown that we can establish a blocking flow in a layered network in $O(nm)$ time.

When a blocking flow x has been established in a network, Dinic's algorithm recomputes the exact distance labels, forms a new layered network, and repeats these computations. The algorithm terminates when as it is forming the new layered networks, it finds that the source is not connected to the sink. It is possible to show that every time Dinic's algorithm forms a new layered network, the distance label of the source node strictly increases. Consequently, Dinic's algorithm forms at most n layered networks and runs in $O(n^2m)$ time.

We now show how to view Dinic's algorithm as a somewhat modified version of the shortest augmenting path algorithm. We make the following three modifications to the shortest augmenting path algorithm.

Modification 1. In operation $\text{retreat}(i)$, we do not change the distance label of node i , but subsequently term node i as *blocked*. A blocked node has no admissible path to the sink node.

Modification 2. We define an arc (i, j) to be admissible if $d(i) = d(j) + 1$, $r_{ij} > 0$, and node j is not blocked.

Modification 3. When the source node is blocked, by performing a backward breadth-first search we recompute the distance labels of all nodes exactly.

We term the computations within two successive recomputations of distance labels as occurring within a single *phase*. We note the following facts about the modified shortest augmenting path algorithm.

1. At the beginning of a phase, when the algorithm recomputes the distance labels $d(\cdot)$, the set of admissible arcs defines a layered network.
2. Each arc (i, j) in the admissible path satisfies $d(i) = d(j) + 1$; therefore, arc (i, j) joins two successive layers of the layered network. As a result, every admissible path is an augmenting path in the layered network.
3. Since we do not update distance labels within a phase, every admissible path has length equal to $d(s)$.

4. The algorithm performs at most m augmentations within a phase because each augmentation causes at least one arc to become inadmissible by reducing its residual capacity to zero, and the algorithm does not create new admissible arcs.
5. A phase ends when the network contains no admissible path from node s to node t . Hence, when the algorithm recomputes distance labels at the beginning of the next phase, $d(s)$ must increase (why?).

The preceding facts show that the modified shortest augmenting path algorithm essentially reduces to Dinic's algorithm. They also show that the distance labels are sufficiently powerful to simulate layered networks. Further, they are simpler to understand than layered networks, easier to manipulate, and lead to more efficient algorithms in practice. Distance labels are also attractive because they are generic solution approaches that find applications in several different algorithms; for example, the generic preflow-push algorithm described next uses distance labels, as does many of its variants described later.

7.6 GENERIC PREFLOW-PUSH ALGORITHM

We now study a class of algorithms, known as *preflow-push* algorithms, for solving the maximum flow problem. These algorithms are more general, more powerful, and more flexible than augmenting path algorithms. The best preflow-push algorithms currently outperform the best augmenting path algorithms in theory as well as in practice. In this section we study the generic preflow-push algorithm. In the following sections we describe special implementations of the generic approach with improved worst-case complexity.

The inherent drawback of the augmenting path algorithms is the computationally expensive operation of sending flow along a path, which requires $O(n)$ time in the worst case. Preflow-push algorithms do not suffer from this drawback and obtain dramatic improvements in the worst-case complexity. To understand this point better, consider the (artificially extreme) example shown in Figure 7.10. When applied to this problem, any augmenting path algorithm would discover 10 augmenting paths, each of length 10, and would augment 1 unit of flow along each of these paths. Observe, however, that although all of these paths share the same first eight arcs, each augmentation traverses all of these arcs. If we could have sent 10 units of flow from node 1 to node 9, and then sent 1 unit of flow along 10 different paths of length 2, we would have saved the repetitive computations in traversing the common set of arcs. This is the essential idea underlying the preflow-push algorithms.

Augmenting path algorithms send flow by augmenting along a path. This basic operation further decomposes into the more elementary operation of sending flow along individual arcs. Thus sending a flow of δ units along a path of k arcs decomposes into k basic operations of sending a flow of δ units along each of the arcs of the path. We shall refer to each of these basic operations as a *push*. The preflow-push algorithms push flows on individual arcs instead of augmenting paths.

Because the preflow-push algorithms push flows along the individual arcs, these algorithms do not satisfy the mass balance constraints (6.1b) at intermediate stages. In fact, these algorithms permit the flow entering a node to exceed the flow leaving

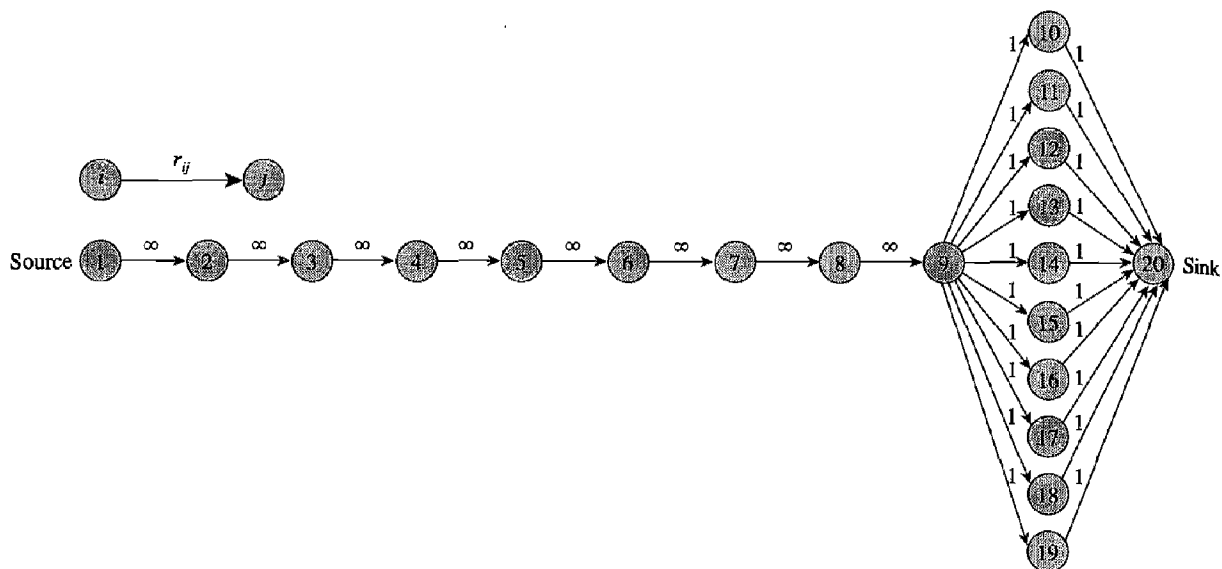


Figure 7.10 Drawback of the augmenting path algorithm.

the node. We refer to any such solution as a *preflow*. More formally, a *preflow* is a function $x: A \rightarrow \mathbf{R}$ that satisfies the flow bound constraint (6.1c) and the following relaxation of (6.1b).

$$\sum_{\{j: (j,i) \in A\}} x_{ji} - \sum_{\{j: (i,j) \in A\}} x_{ij} \geq 0 \quad \text{for all } i \in N - \{s, t\}.$$

The preflow-push algorithms maintain a preflow at each intermediate stage. For a given preflow x , we define the *excess* of each node $i \in N$ as

$$e(i) = \sum_{\{j: (j,i) \in A\}} x_{ji} - \sum_{\{j: (i,j) \in A\}} x_{ij}.$$

In a preflow, $e(i) \geq 0$ for each $i \in N - \{s, t\}$. Moreover, because no arc emanates from node t in the preflow push algorithms, $e(t) \geq 0$ as well. Therefore node s is the only node with negative excess.

We refer to a node with a (strictly) positive excess as an *active* node and adopt the convention that the source and sink nodes are never active. The augmenting path algorithms always maintain feasibility of the solution and strive toward optimality. In contrast, preflow-push algorithms strive to achieve feasibility. In a preflow-push algorithm, the presence of active nodes indicates that the solution is infeasible. Consequently, the basic operation in this algorithm is to select an active node and try to remove its excess by pushing flow to its neighbors. But to which nodes should the flow be sent? Since ultimately we want to send flow to the sink, we push flow to the nodes that are *closer* to sink. As in the shortest augmenting path algorithm, we measure closeness with respect to the current distance labels, so sending flow closer to the sink is equivalent to pushing flow on admissible arcs. Thus we send flow only on admissible arcs. If the active node we are currently considering has no admissible arc, we increase its distance label so that we create

at least one admissible arc. The algorithm terminates when the network contains no active node. The preflow-push algorithm uses the subroutines shown in Figure 7.11.

```

procedure preprocess;
begin
   $x := 0$ ;
  compute the exact distance labels  $d(i)$ ;
   $x_{sj} := u_{sj}$  for each arc  $(s, j) \in A(s)$ ;
   $d(s) := n$ ;
end;

(a)

procedure push/relabel(i);
begin
  if the network contains an admissible arc  $(i, j)$  then
    push  $\delta := \min\{e(i), r_{ij}\}$  units of flow from node  $i$  to node  $j$ 
  else replace  $d(i)$  by  $\min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\}$ ;
end;

(b)

```

Figure 7.11 Subroutines of the preflow-push algorithm.

A push of δ units from node i to node j decreases both $e(i)$ and r_{ij} by δ units and increases both $e(j)$ and r_{ji} by δ units. We say that a push of δ units of flow on an arc (i, j) is *saturating* if $\delta = r_{ij}$ and is *nonsaturating* otherwise. A nonsaturating push at node i reduces $e(i)$ to zero. We refer to the process of increasing the distance label of a node as a *relabel* operation. The purpose of the relabel operation is to create at least one admissible arc on which the algorithm can perform further pushes.

The generic version of the preflow-push algorithm (Figure 7.12) combines the subroutines just described.

```

algorithm preflow-push;
begin
  preprocess;
  while the network contains an active node do
    begin
      select an active node  $i$ ;
      push/relabel(i);
    end;
  end;

```

Figure 7.12 Generic preflow-push algorithm.

It might be instructive to visualize the generic preflow-push algorithm in terms of a physical network: Arcs represent flexible water pipes, nodes represent joints, and the distance function measures how far nodes are above the ground. In this network we wish to send water from the source to the sink. In addition, we visualize flow in an admissible arc as water flowing downhill. Initially, we move the source node upward, and water flows to its neighbors. In general, water flows downhill towards the sink; however, occasionally flow becomes trapped locally at a node that has no downhill neighbors. At this point we move the node upward, and again water

flows downhill toward the sink. Eventually, no more flow can reach the sink. As we continue to move nodes upward, the remaining excess flow eventually flows back toward the source. The algorithm terminates when all the water flows either into the sink or flows back to the source.

The preprocessing operation accomplishes several important tasks. First, it gives each node adjacent to node s a positive excess, so that the algorithm can begin by selecting some node with a positive excess. Second, since the preprocessing operation saturates all the arcs incident to node s , none of these arcs is admissible and setting $d(s) = n$ will satisfy the validity condition (7.2). Third, since $d(s) = n$, Property 7.2 implies that the residual network contains no directed path from node s to node t . Since distance labels are nondecreasing, we also guarantee that in subsequent iterations the residual network will never contain a directed path from node s to node t , and so we will never need to push flow from node s again.

To illustrate the generic preflow-push algorithm, consider the example given in Figure 7.13(a). Figure 7.13(b) specifies the preflow determined by the preprocess operation.

Iteration 1. Suppose that the algorithm selects node 2 for the push/relabel operation. Arc (2, 4) is the only admissible arc and the algorithm performs a push of value $\delta = \min\{e(2), r_{24}\} = \min\{2, 1\} = 1$. This push is saturating. Figure 7.13(c) gives the residual network at this stage.

Iteration 2. Suppose that the algorithm again selects node 2. Since no admissible arc emanates from node 2, the algorithm performs a relabel operation and gives node 2 a new distance label $d(2) = \min\{d(3) + 1, d(1) + 1\} = \min\{2, 5\} = 2$. The new residual network is the same as the one shown in Figure 7.13(c) except that $d(2) = 2$ instead of 1.

Iteration 3. Suppose that this time the algorithm selects node 3. Arc (3, 4) is the only admissible arc emanating from node 3, the algorithm performs a push of value $\delta = \min\{e(3), r_{34}\} = \min\{4, 5\} = 4$. This push is nonsaturating. Figure 7.13(d) gives the residual network at the end of this iteration.

Iteration 4. The algorithm selects node 2 and performs a nonsaturating push of value $\delta = \min\{1, 3\} = 1$, obtaining the residual network given in Figure 7.13(e).

Iteration 5. The algorithm selects node 3 and performs a saturating push of value $\delta = \min\{1, 1\} = 1$ on arc (3, 4), obtaining the residual network given in Figure 7.13(f).

Now the network contains no active node and the algorithm terminates. The maximum flow value in the network is $e(4) = 6$.

Assuming that the generic preflow-push algorithm terminates, we can easily show that it finds a maximum flow. The algorithm terminates when the excess resides at the source or at the sink, implying that the current preflow is a flow. Since $d(s) = n$, the residual network contains no path from the source to the sink. This condition is the termination criterion of the augmenting path algorithm, and the excess residing at the sink is the maximum flow value.

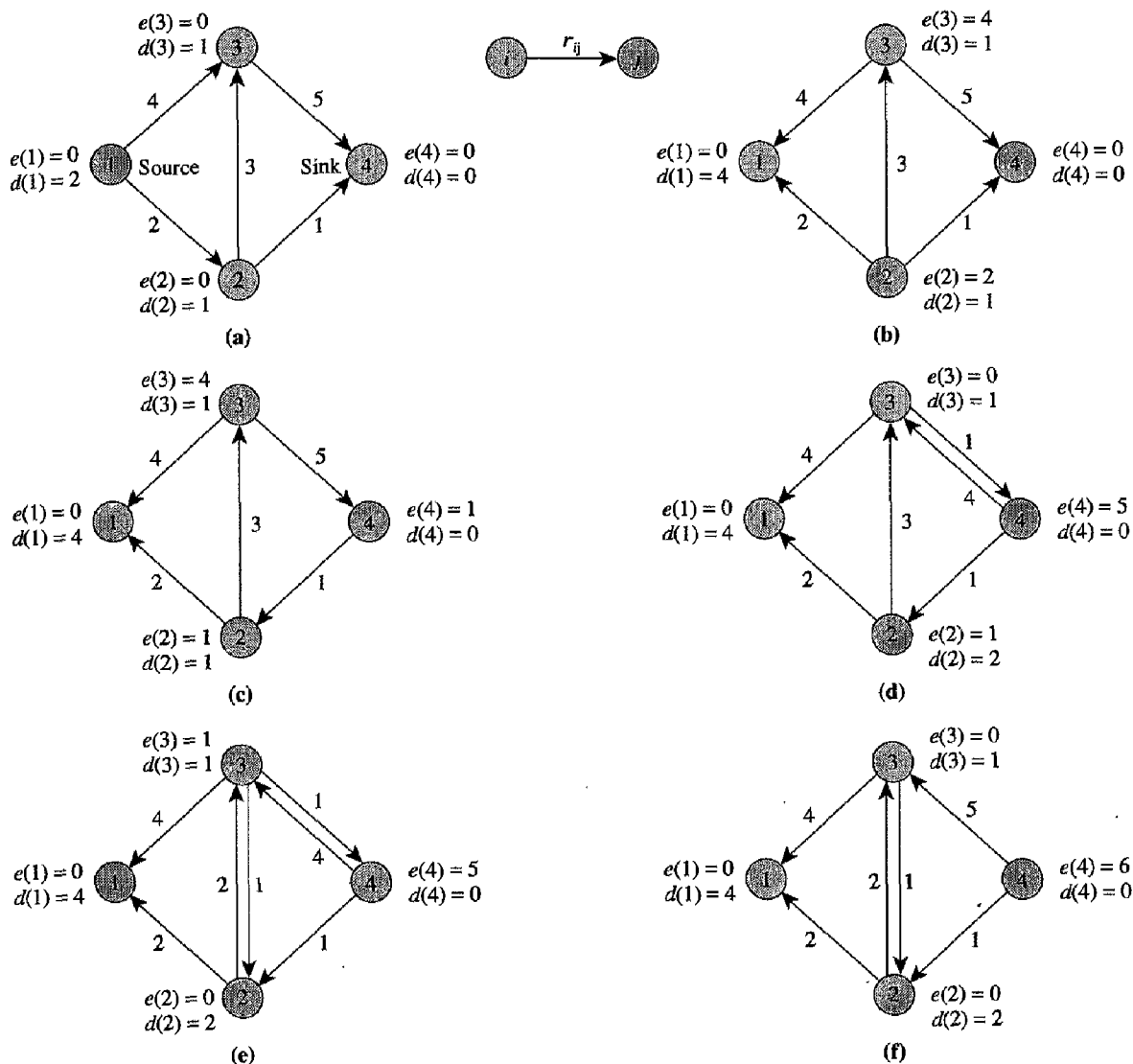


Figure 7.13 Illustrating the generic preflow-push algorithm.

Complexity of the Algorithm

To analyze the complexity of the algorithm, we begin by establishing one important result: distance labels are always valid and do not increase “too many” times. The first of these conclusions follows from Lemma 7.5, because, as in the shortest augmenting path algorithm, the preflow-push algorithm pushes flow only on admissible arcs and relabels a node only when no admissible arc emanates from it. The second conclusion follows from the following lemma.

Lemma 7.11. *At any stage of the preflow-push algorithm, each node i with positive excess is connected to node s by a directed path from node i to node s in the residual network.*

Proof. Notice that for a preflow x , $e(s) \leq 0$ and $e(i) \geq 0$ for all $i \in N - \{s\}$. By the flow decomposition theorem (see Theorem 3.5), we can decompose any preflow x with respect to the original network G into nonnegative flows along (1) paths from node s to node t , (2) paths from node s to active nodes, and (3) flows around directed cycles. Let i be an active node relative to the preflow x in G . The flow decomposition of x must contain a path P from node s to node i , since the paths from node s to node t and the flows around cycles do not contribute to the excess at node i . The residual network contains the reversal of P (P with the orientation of each arc reversed), so a directed path from node i to node s . ♦

This lemma implies that during a relabel operation, the algorithm does not minimize over an empty set.

Lemma 7.12. *For each node $i \in N$, $d(i) < 2n$.*

Proof. The last time the algorithm relabeled node i , the node had a positive excess, so the residual network contained a path P of length at most $n - 2$ from node i to node s . The fact that $d(s) = n$ and that $d(k) \leq d(l) + 1$ for every arc (k, l) in the path P implies that $d(i) \leq d(s) + |P| < 2n$. ♦

Since each time the algorithm relabels node i , $d(i)$ increases by at least 1 unit, we have established the following result.

Lemma 7.13. *Each distance label increases at most $2n$ times. Consequently, the total number of relabel operations is at most $2n^2$.* ♦

Lemma 7.14. *The algorithm performs at most nm saturating pushes.*

Proof. This result follows directly from Lemmas 7.12 and 7.8. ♦

In view of Lemma 7.7, Lemma 7.13 implies that the total time needed to identify admissible arcs and to perform relabel operations is $O(nm)$. We next count the number of nonsaturating pushes performed by the algorithm.

Lemma 7.15. *The generic preflow-push algorithm performs $O(n^2m)$ nonsaturating pushes.*

Proof. We prove the lemma using an argument based on potential functions (see Section 3.2). Let I denote the set of active nodes. Consider the potential function $\Phi = \sum_{i \in I} d(i)$. Since $|I| < n$, and $d(i) < 2n$ for all $i \in I$, the initial value of Φ (after the preprocess operation) is at most $2n^2$. At the termination of the algorithm, Φ is zero. During the push/relabel(i) operation, one of the following two cases must apply:

Case 1. The algorithm is unable to find an admissible arc along which it can push flow. In this case the distance label of node i increases by $\epsilon \geq 1$ units. This operation increases Φ by at most ϵ units. Since the total increase in $d(i)$ for each node i throughout the execution of the algorithm is bounded by $2n$, the total increase in Φ due to increases in distance labels is bounded by $2n^2$.

Case 2. The algorithm is able to identify an arc on which it can push flow, so it performs a saturating push or a nonsaturating push. A saturating push on arc (i, j) might create a new excess at node j , thereby increasing the number of active nodes by 1, and increasing Φ by $d(j)$, which could be as much as $2n$ per saturating push, and so $2n^2m$ over all saturating pushes. Next note that a nonsaturating push on arc (i, j) does not increase $|I|$. The nonsaturating push will decrease Φ by $d(i)$ since i becomes inactive, but it simultaneously increases Φ by $d(j) = d(i) - 1$ if the push causes node j to become active, the total decrease in Φ being of value 1. If node j was active before the push, Φ decreases by an amount $d(i)$. Consequently, net decrease in Φ is at least 1 unit per nonsaturating push.

We summarize these facts. The initial value of Φ is at most $2n^2$ and the maximum possible increase in Φ is $2n^2 + 2n^2m$. Each nonsaturating push decreases Φ by at least 1 unit and Φ always remains nonnegative. Consequently, the algorithm can perform at most $2n^2 + 2n^2 + 2n^2m = O(n^2m)$ nonsaturating pushes, proving the lemma. ♦

Finally, we indicate how the algorithm keeps track of active nodes for the push/relabel operations. The algorithm maintains a set LIST of active nodes. It adds to LIST those nodes that become active following a push and are not already in LIST, and deletes from LIST nodes that become inactive following a nonsaturating push. Several data structures (e.g., doubly linked lists) are available for storing LIST so that the algorithm can add, delete, or select elements from it in $O(1)$ time. Consequently, it is easy to implement the preflow-push algorithm in $O(n^2m)$ time. We have thus established the following theorem.

Theorem 7.16. *The generic preflow-push algorithm runs in $O(n^2m)$ time.* ♦

Several modifications to the generic preflow-push algorithm might improve its empirical performance. We define a maximum preflow as a preflow with the maximum possible flow into the sink. As stated, the generic preflow-push algorithm performs push/relabel operations at active nodes until all the excess reaches the sink node or returns to the source node. Typically, the algorithm establishes a maximum preflow long before it establishes a maximum flow; the subsequent push/relabel operations increase the distance labels of the active nodes until they are sufficiently higher than n so they can push their excesses back to the source node (whose distance label is n). One possible modification in the preflow-push algorithm is to maintain a set N' of nodes that satisfy the property that the residual network contains no path from a node in N' to the sink node t . Initially, $N' = \{s\}$ and, subsequently, whenever the distance label of a node is greater than or equal to n , we add it to N' . Further, we do not perform push/relabel operations for nodes in N' and terminate the algorithm when all nodes in $N - N'$ are inactive. At termination, the current preflow x is also an optimal preflow. At this point we convert the maximum preflow x into a maximum flow using any of the methods described in Exercise 7.11. Empirical tests have found that this two-phase approach often substantially reduces the running times of preflow push algorithms.

One sufficient condition for adding a node j to N' is $d(j) \geq n$. Unfortunately,

this simple approach is not very effective and does not substantially reduce the running time of the algorithm. Another approach is to occasionally perform a reverse breadth-first search of the residual network to obtain exact distance labels and add all those nodes to N' that do not have any directed path to the sink. Performing this search occasionally, that is, after αn relabel operations for some constant α , does not effect the worst-case complexity of the preflow-push algorithm (why?) but improves the empirical behavior of the algorithm substantially.

A third approach is to let $\text{numb}(k)$ denote the number of nodes whose distance label is k . As discussed in Section 7.4, we can update the array $\text{numb}(\cdot)$ in $O(1)$ steps per relabel operation. Moreover, whenever $\text{numb}(k') = 0$ for some k' , any node j with $d(j) > k'$ is disconnected from the set of nodes i with $d(i) < k'$ in the residual network. At this point, we can increase the distance labels of each of these nodes to n and the distance labels will still be valid (why?). Equivalently, we can add any node j with $d(j) > k'$ to the set N' . The array $\text{numb}(\cdot)$ is easy to implement, and its use is quite effective in practice.

Specific Implementations of Generic Preflow-Push Algorithm

The running time of the generic preflow-push algorithm is comparable to the bound of the shortest augmenting path algorithm. However, the preflow-push algorithm has several nice features: in particular, its flexibility and its potential for further improvements. By specifying different rules for selecting active nodes for the push/relabel operations, we can derive many different algorithms, each with different worst-case complexity than the generic version of the algorithm. The bottleneck operation in the generic preflow-push algorithm is the number of nonsaturating pushes and many specific rules for examining active nodes can produce substantial reductions in the number of nonsaturating pushes. We consider the following three implementations.

1. *FIFO preflow-push algorithm.* This algorithm examines the active nodes in the first-in, first-out (FIFO) order. We shall show that this algorithm runs in $O(n^3)$ time.
2. *Highest-label preflow-push algorithm.* This algorithm always pushes from an active node with the highest value of the distance label. We shall show that this algorithm runs in $O(n^2 m^{1/2})$ time. Observe that this time is better than $O(n^3)$ for all problem densities.
3. *Excess scaling algorithm.* This algorithm pushes flow from a node with sufficiently large excess to a node with sufficiently small excess. We shall show that the excess scaling algorithm runs in $O(nm + n^2 \log U)$ time. For problems that satisfy the similarity assumption (see Section 3.2), this time bound is better than that of the two preceding algorithms.

We might note that the time bounds for all these preflow-push algorithms are tight (except the excess scaling algorithm); that is, for some classes of networks the generic preflow-push algorithm, the FIFO algorithm, and the highest-label preflow-push algorithms do perform as many computations as indicated by their worst-case

time bounds. These examples show that we cannot improve the time bounds of these algorithms by a more clever analysis.

7.7 FIFO PREFLOW-PUSH ALGORITHM

Before we describe the FIFO implementation of the preflow-push algorithm, we define the concept of a *node examination*. In an iteration, the generic preflow-push algorithm selects a node, say node i , and performs a saturating push or a nonsaturating push, or relabels the node. If the algorithm performs a saturating push, then node i might still be active, but it is not mandatory for the algorithm to select this node again in the next iteration. The algorithm might select another node for the next push/relabel operation. However, it is easy to incorporate the rule that whenever the algorithm selects an active node, it keeps pushing flow from that node until either the node's excess becomes zero or the algorithm relabels the node. Consequently, the algorithm might perform several saturating pushes followed either by a non-saturating push or a relabel operation. We refer to this sequence of operations as a node examination. We shall henceforth assume that every preflow-push algorithm adopts this rule for selecting nodes for the push/relabel operation.

The FIFO preflow-push algorithm examines active nodes in the FIFO order. The algorithm maintains the set LIST as a queue. It selects a node i from the front of LIST, performs pushes from this node, and adds newly active nodes to the rear of LIST. The algorithm examines node i until either it becomes inactive or it is relabeled. In the latter case, we add node i to the rear of the queue. The algorithm terminates when the queue of active nodes is empty.

We illustrate the FIFO preflow-push algorithm using the example shown in Figure 7.14(a). The preprocess operation creates an excess of 10 units at each of

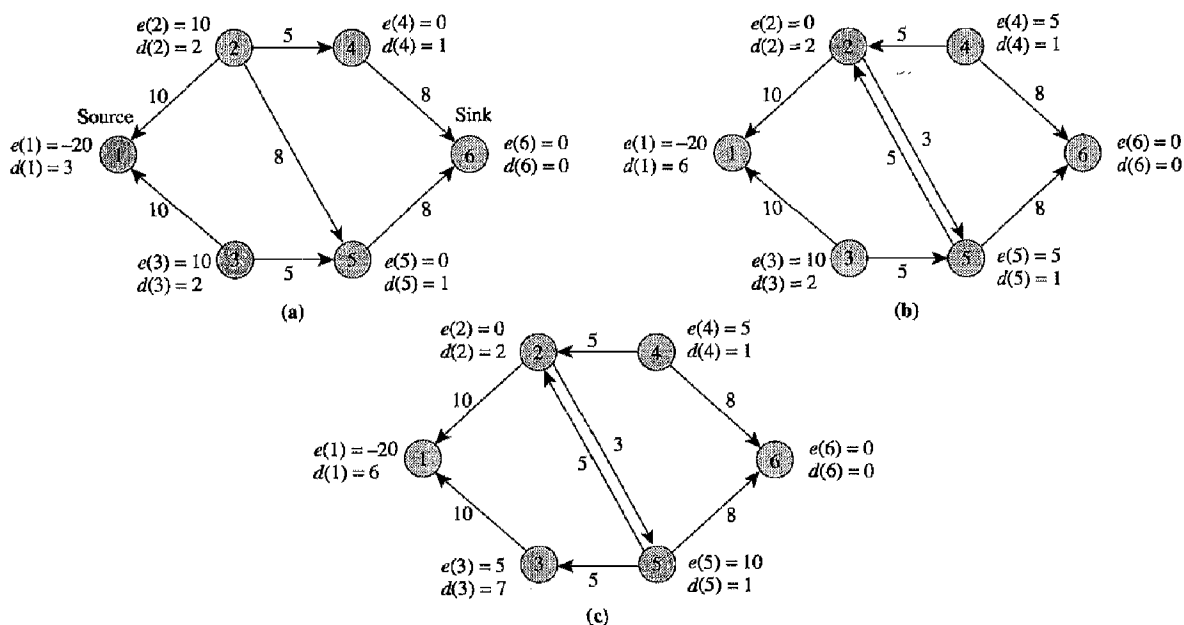


Figure 7.14 Illustrating the FIFO preflow-push algorithm.

the nodes 2 and 3. Suppose that the queue of active nodes at this stage is $LIST = \{2, 3\}$. The algorithm removes node 2 from the queue and examines it. Suppose that it performs a saturating push of 5 units on arc $(2, 4)$ and a nonsaturating push of 5 units on arc $(2, 5)$ [see Figure 7.14(b)]. As a result of these pushes, nodes 4 and 5 become active and we add these nodes to the queue in this order, obtaining $LIST = \{3, 4, 5\}$. The algorithm next removes node 3 from the queue. While examining node 3, the algorithm performs a saturating push of 5 units on arc $(3, 5)$, followed by a relabel operation of node 3 [see Figure 7.14(c)]. The algorithm adds node 3 to the queue, obtaining $LIST = \{4, 5, 3\}$. We encourage the reader to complete the solution of this example.

To analyze the worst-case complexity of the FIFO preflow-push algorithm, we partition the total number of node examinations into different phases. The first phase consists of node examinations for those nodes that become active during the pre-process operation. The second phase consists of the node examinations of all the nodes that are in the queue after the algorithm has examined the nodes in the first phase. Similarly, the third phase consists of the node examinations of all the nodes that are in the queue after the algorithm has examined the nodes in the second phase, and so on. For example, in the preceding illustration, the first phase consists of the node examinations of the set $\{2, 3\}$, and the second phase consists of the node examinations of the set $\{4, 5, 3\}$. Observe that the algorithm examines any node at most once during a phase.

We will now show that the algorithm performs at most $2n^2 + n$ phases. Each phase examines any node at most once and each node examination performs at most one nonsaturating push. Therefore, a bound of $2n^2 + n$ on the total number of phases would imply a bound of $O(n^3)$ on the number of nonsaturating pushes. This result would also imply that the FIFO preflow-push algorithm runs in $O(n^3)$ time because the bottleneck operation in the generic preflow-push algorithm is the number of nonsaturating pushes.

To bound the number of phases in the algorithm, we consider the total change in the potential function $\Phi = \max\{d(i) : i \text{ is active}\}$ over an entire phase. By the “total change” we mean the difference between the initial and final values of the potential function during a phase. We consider two cases.

Case 1. The algorithm performs at least one relabel operation during a phase. Then Φ might increase by as much as the maximum increase in any distance label. Lemma 7.13 implies that the total increase in Φ over all the phases is at most $2n^2$.

Case 2. The algorithm performs no relabel operation during a phase. In this case the excess of every node that was active at the beginning of the phase moves to nodes with smaller distance labels. Consequently, Φ decreases by at least 1 unit.

Combining Cases 1 and 2, we find that the total number of phases is at most $2n^2 + n$; the second term corresponds to the initial value of Φ , which could be at most n . We have thus proved the following theorem.

Theorem 7.17. *The FIFO preflow-push algorithm runs in $O(n^3)$ time. ♦*

7.8 HIGHEST-LABEL PREFLOW-PUSH ALGORITHM

The highest-label preflow-push algorithm always pushes flow from an active node with the highest distance label. It is easy to develop an $O(n^3)$ bound on the number of nonsaturating pushes for this algorithm. Let $h^* = \max\{d(i) : i \text{ is active}\}$. The algorithm first examines nodes with distance labels equal to h^* and pushes flow to nodes with distance labels equal to $h^* - 1$, and these nodes, in turn, push flow to nodes with distance labels equal to $h^* - 2$, and so on, until either the algorithm relabels a node or it has exhausted all the active nodes. When it has relabeled a node, the algorithm repeats the same process. Note that if the algorithm does not relabel any node during n consecutive node examinations, all the excess reaches the sink (or the source) and the algorithm terminates. Since the algorithm performs at most $2n^2$ relabel operations (by Lemma 7.13), we immediately obtain a bound of $O(n^3)$ on the number of node examinations. Since each node examination entails at most one nonsaturating push, the highest-label preflow-push algorithm performs $O(n^3)$ nonsaturating pushes. (In Exercise 7.20 we consider a potential function argument that gives the same bound on the number of nonsaturating pushes.)

The preceding discussion is missing one important detail: How do we select a node with the highest distance label without expending too much effort? We use the following data structure. For each $k = 1, 2, \dots, 2n - 1$, we maintain the list

$$\text{LIST}(k) = \{i : i \text{ is active and } d(i) = k\},$$

in the form of either linked stacks or linked queues (see Appendix A). We define a variable *level* that is an upper bound on the highest value of k for which $\text{LIST}(k)$ is nonempty. To determine a node with the highest distance label, we examine the lists $\text{LIST}(\text{level})$, $\text{LIST}(\text{level}-1)$, \dots , until we find a nonempty list, say $\text{LIST}(p)$. We set *level* equal to p and select any node in $\text{LIST}(p)$. Moreover, if the distance label of a node increases while the algorithm is examining it, we set *level* equal to the new distance label of the node. Observe that the total increase in *level* is at most $2n^2$ (from Lemma 7.13), so the total decrease is at most $2n^2 + n$. Consequently, scanning the lists $\text{LIST}(\text{level})$, $\text{LIST}(\text{level}-1)$, \dots , in order to find the first nonempty list is not a bottleneck operation.

The highest-label preflow-push algorithm is currently the most efficient method for solving the maximum flow problem in practice because it performs the least number of nonsaturating pushes. To illustrate intuitively why the algorithm performs so well in practice, we consider the maximum flow problem given in Figure 7.15(a). The preprocess operation creates an excess of 1 unit at each node $2, 3, \dots, n - 1$ [see Figure 7.15(b)]. The highest-label preflow-push algorithm examines nodes $2, 3, \dots, n - 1$, in this order and pushes all the excess to the sink node. In contrast, the FIFO preflow-push algorithm might perform many more pushes. Suppose that at the end of the preprocess operation, the queue of active nodes is $\text{LIST} = \{n - 1, n - 2, \dots, 3, 2\}$. Then the algorithm would examine each of these nodes in the first phase and would obtain the solution depicted in Figure 7.15(c). At this point, $\text{LIST} = \{n - 1, n - 2, \dots, 4, 3\}$. It is easy to show that overall the algorithm would perform $n - 2$ phases and use $(n - 2) + (n - 3) + \dots + 1 = \Omega(n^2)$ nonsaturating pushes.

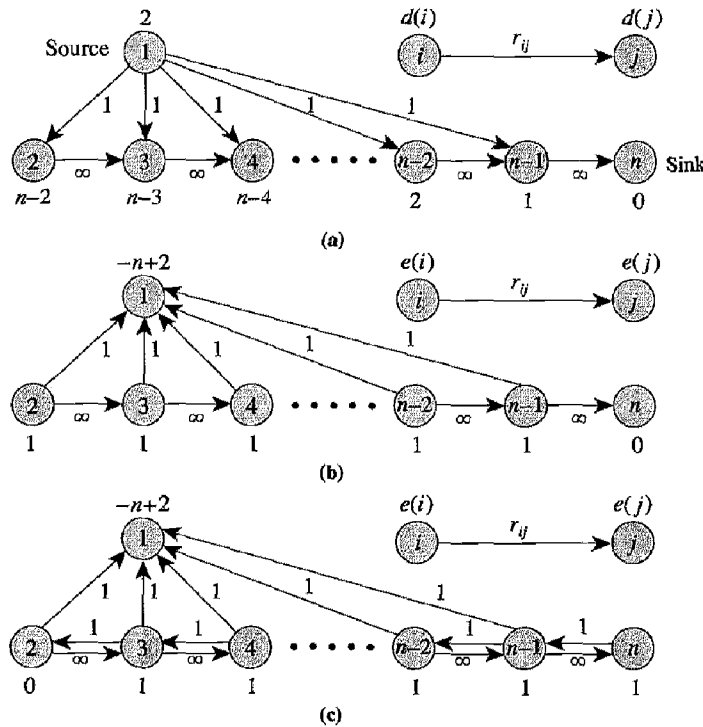


Figure 7.15 Bad example for the FIFO preflow-push algorithm: (a) initial residual network; (b) network after the preprocess operation; (c) network after one phase of the FIFO preflow-push algorithm.

Although the preceding example is rather extreme, it does illustrate the advantage in pushing flows from active nodes with the highest distance label. In our example the FIFO algorithm selects an excess and pushes it all the way to the sink. Then it selects another excess and pushes it to the sink, and repeats this process until no node contains any more excess. On the other hand, the highest-label preflow-push algorithm starts at the highest level and pushes all the excess at this level to the next lower level and repeats this process. As the algorithm examines nodes with lower and lower distance labels, it accumulates the excesses and pushes this accumulated excess toward the sink. Consequently, the highest-label preflow-push algorithm avoids repetitive pushes on arcs carrying a small amount of flow.

This nice feature of the highest-label preflow-push algorithm also translates into a tighter bound on the number of nonsaturating pushes. The bound of $O(n^3)$ on the number of nonsaturating pushes performed by the algorithm is rather loose and can be improved by a more clever analysis. We now show that the algorithm in fact performs $O(n^2 m^{1/2})$ nonsaturating pushes. The proof of this theorem is somewhat complex and the reader can skip it without any loss of continuity.

At every state of the preflow-push algorithm, each node other than the sink has at most one current arc which, by definition, must be admissible. We denote this collection of current arcs by the set F . The set F has at most $n - 1$ arcs, has at most one outgoing arc per node, and does not contain any cycle (why?). These results imply that F defines a *forest*, which we subsequently refer to as the *current forest*. Figure 7.16 gives an example of a current forest. Notice that each tree in the forest is a rooted tree, the root being a node with no outgoing arc.

Before continuing, let us introduce some additional notation. For any node $i \in N$, we let $D(i)$ denote the set of descendants of that node in F (we refer the

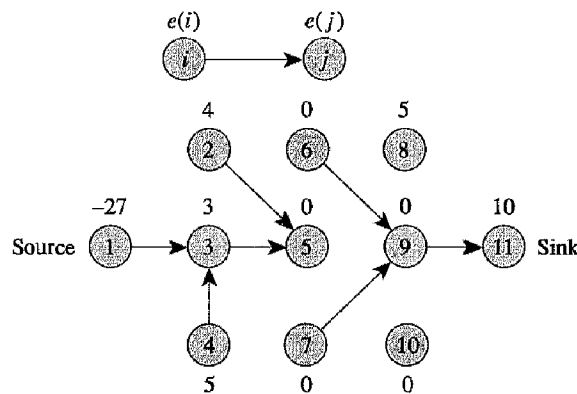


Figure 7.16 Example of a current-forest.

reader to Section 2.2 for the definition of descendants in a rooted tree). For example, in Figure 7.16, $D(1) = \{1\}$, $D(2) = \{2\}$, $D(3) = \{1, 3, 4\}$, $D(4) = \{4\}$, and $D(5) = \{1, 2, 3, 4, 5\}$. Notice that distance label of the descendants of any node i will be higher than $d(i)$. We refer to an active node with no active descendants (other than itself) as a *maximal active node*. In Figure 7.16 the nodes 2, 4, and 8 are the only maximal active nodes. Let H denote the set of maximal active nodes. Notice that two maximal active nodes have distinct descendants. Also notice that the highest-label preflow-push algorithm always pushes flow from a maximal active node.

We obtain the time bound of $O(n^2 m^{1/2})$ for the highest-label preflow-push algorithm using a potential function argument. The argument relies on a parameter K , whose optimal value we select later. Our potential function is $\Phi = \sum_{i \in H} \Phi(i)$, with $\Phi(i)$ defined as $\Phi(i) = \max\{0, K + 1 - |D(i)|\}$. Observe that for any node i , $\Phi(i)$ is at most K [because $|D(i)| \geq 1$]. Also observe that Φ changes whenever the set H of maximal active nodes changes or $|D(i)|$ changes for a maximal active node i .

We now study the effect of various operations performed by the preflow-push algorithm on the potential function Φ . As the algorithm proceeds, it changes the set of current arcs, performs saturating and nonsaturating pushes, and relabels nodes. All these operations have an effect on the value of Φ . By observing the consequence of all these operations on Φ , we will obtain a bound on the number of nonsaturating pushes.

First, consider a nonsaturating push on an arc (i, j) emanating from a maximal active node i . Notice that a nonsaturating push takes place on a current arc and does not change the current forest; it simply moves the excess from node i to node j [see Figure 7.17(a) for a nonsaturating push on the arc $(3, 4)$]. As a result of the push, node i becomes inactive and node j might become a new maximal active node. Since $|D(j)| > |D(i)|$, this push decreases $\Phi(i) + \Phi(j)$ by at least 1 unit if $|D(i)| \leq K$ and does not change $\Phi(i) + \Phi(j)$ otherwise.

Now consider a saturating push on the arc (i, j) emanating from a maximal active node i . As a result of the push, arc (i, j) becomes inadmissible and drops out of the current forest [see Figure 7.17(b) for a saturating push on the arc $(1, 3)$]. Node i remains a maximal active node and node j might also become a maximal active node. Consequently, this operation might increase Φ by upto K units.

Next consider the relabeling of a maximal active node i . We relabel a node when it has no admissible arc; therefore, no current arc emanates from this node.

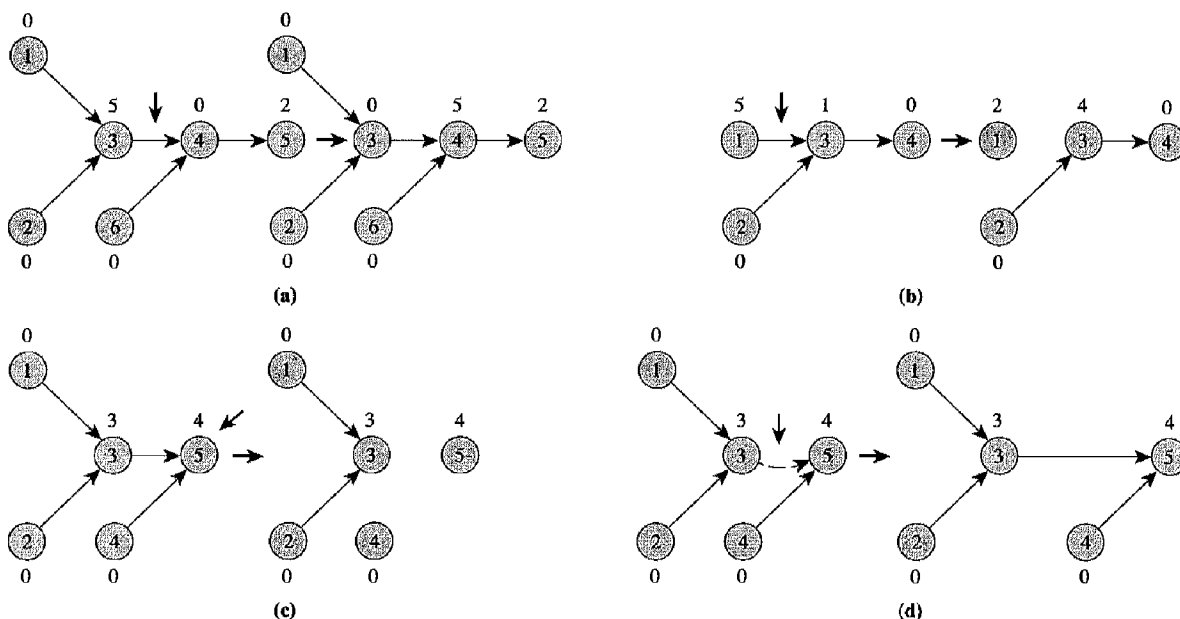


Figure 7.17 (a) Nonsaturating push on arc $(d, 4)$; (b) saturating push on arc $(1, 3)$; (c) relabel of node 5; (d) addition of the arc $(3, 5)$ to the forest.

As a consequence, node i must be a root node in the current forest. Moreover, since node i is a maximal active node, none of its proper descendants can be active. After the algorithm has relabeled node i , all incoming arcs at node i become inadmissible; therefore, all the current arcs entering node i will no longer belong to the current forest [see Figure 7.17(c)]. Clearly, this change cannot create any new maximal active nodes. The relabel operation, however, decreases the number of descendants of node i to one. Consequently, $\Phi(i)$ can increase by at most K .

Finally, consider the introduction of new current arcs in the current forest. The addition of new arcs to the forest does not create any new maximal active nodes. It might, in fact, remove some maximal active nodes and increase the number of descendants of some nodes [see Figure 7.17(d)]. In both cases the potential Φ does not increase. We summarize the preceding discussion in the form of the following property.

Property 7.18

- (a) A nonsaturating push from a maximal active node i does not increase Φ ; it decreases Φ by at least 1 unit if $|D(i)| \leq K$.
- (b) A saturating push from a maximal active node i can increase Φ by at most K units.
- (c) The relabeling of a maximal active node i can increase Φ by at most K units.
- (d) Introducing current arcs does not increase Φ .

For the purpose of worst-case analysis, we define the concept of phases. A *phase* consists of the sequence of pushes between two consecutive relabel operations. Lemma 7.13 implies that the algorithm contains $O(n^2)$ phases. We call a phase *cheap* if it performs at most $2n/K$ nonsaturating pushes, and *expensive* otherwise.

Clearly, the number of nonsaturating pushes in cheap phases is at most $O(n^2 \cdot 2n/K) = O(n^3/K)$. To obtain a bound on the nonsaturating pushes in expensive phases, we use an argument based on the potential function Φ .

By definition, an expensive phase performs at least $2n/K$ nonsaturating pushes. Since the network can contain at most n/K nodes with K descendants or more, at least n/K nonsaturating pushes must be from nodes with fewer than K descendants. The highest-label preflow-push algorithm always performs push/relabel operation on a maximal active node; consequently, Property 7.18 applies. Property 7.18(a) implies that each of these nonsaturating pushes produces a decrease in Φ of at least 1. So Properties 7.18(b) and (c) imply that the total increase in Φ due to saturating pushes and relabels is at most $O(nmK)$. Therefore, the algorithm can perform $O(nmK)$ nonsaturating pushes in expensive phases.

To summarize this discussion, we note that cheap phases perform $O(n^3/K)$ nonsaturating pushes and expensive phases perform $O(nmK)$ nonsaturating pushes. We obtain the optimal value of K by balancing both terms (see Section 3.2), that is, when both the terms are equal: $n^3/K = nmK$ or $K = n/m^{1/2}$. For this value of K , the number of nonsaturating pushes is $O(n^2m^{1/2})$. We have thus established the following result.

Theorem 7.19. *The highest-label preflow-push algorithm performs $O(n^2m^{1/2})$ nonsaturating pushes and runs in the same time.* ♦

7.9 EXCESS SCALING ALGORITHM

The generic preflow-push algorithm allows flow at each intermediate step to violate the mass balance equations. By pushing flows from active nodes, the algorithm attempts to satisfy the mass balance equations. The function $e_{\max} = \max\{e(i) : i \text{ is an active node}\}$ provides one measure of the infeasibility of a preflow. Note that during the execution of the generic algorithm, we would observe no particular pattern in the values of e_{\max} , except that e_{\max} eventually decreases to value 0. In this section we develop an *excess scaling technique* that systematically reduces the value of e_{\max} to 0.

The excess scaling algorithm is similar to the capacity scaling algorithm we discussed in Section 7.3. Recall that the generic augmenting path algorithm performs $O(nU)$ augmentations and the capacity scaling algorithm reduces this number to $O(m \log U)$ by assuring that each augmentation carries a “sufficiently large” amount of flow. Similarly, in the generic preflow-push algorithm, nonsaturating pushes carrying small amount of flow bottleneck the algorithm in theory. The excess scaling algorithm assures that each nonsaturating push carries a “sufficiently large” amount of flow and so the number of nonsaturating pushes is “sufficiently small.”

Let Δ denote an upper bound on e_{\max} ; we refer to this bound as the *excess dominator*. We refer to a node with $e(i) \geq \Delta/2 \geq e_{\max}/2$ as a node with *large excess*, and as a node with *small excess* otherwise. The excess scaling algorithm always pushes flow from a node with a large excess. This choice assures that during nonsaturating pushes, the algorithm sends relatively large excess closer to the sink.

The excess scaling algorithm also does not allow the maximum excess to increase beyond Δ . This algorithmic strategy might prove to be useful for the following

reason. Suppose that several nodes send flow to a single node j , creating a very large excess. It is likely that node j cannot send the accumulated flow closer to the sink, and thus the algorithm will need to increase its distance label and return much of its excess back to the nodes it came from. Thus pushing too much flow to any node is also likely to be a wasted effort.

The two conditions we have discussed—that each nonsaturating push must carry at least $\Delta/2$ units of flow and that no excess should exceed Δ —imply that we need to select the active nodes for push/relabel operations carefully. The following selection rule is one that assures that we achieve these objectives.

Node Selection Rule. *Among all nodes with a large excess, select a node with the smallest distance label (breaking ties arbitrarily).*

We are now in a position to give, in Figure 7.18, a formal description of the excess scaling algorithm.

The excess scaling algorithm uses the same $\text{push/relabel}(i)$ operation as the generic preflow-push algorithm, but with one slight difference. Instead of pushing $\delta = \min\{e(i), r_{ij}\}$ units of flow, it pushes $\delta = \min\{e(i), r_{ij}, \Delta - e(j)\}$ units. This change ensures that the algorithm permits no excess to exceed Δ .

The algorithm performs a number of scaling phases with the value of the excess dominator Δ decreasing from phase to phase. We refer to a specific scaling phase with a particular value of Δ as a Δ -scaling phase. Initially, $\Delta = 2^{\lceil \log U \rceil}$. Since the logarithm is of base 2, $U \leq \Delta \leq 2U$. During the Δ -scaling phase, $\Delta/2 < e_{\max} \leq \Delta$; the value of e_{\max} might increase or decrease during the phase. When $e_{\max} \leq \Delta/2$, we begin a new scaling phase. After the algorithm has performed $\lceil \log U \rceil + 1$ scaling phases, e_{\max} decreases to value 0 and we obtain the maximum flow.

Lemma 7.20. *The algorithm satisfies the following two conditions:*

- (a) *Each nonsaturating push sends at least $\Delta/2$ units of flow.*
- (b) *No excess ever exceeds Δ .*

Proof. Consider a nonsaturating push on arc (i, j) . Since arc (i, j) is admissible, $d(j) < d(i)$. Moreover, since node i is a node with the smallest distance label among all nodes with a large excess, $e(i) \geq \Delta/2$ and $e(j) < \Delta/2$. Since this push is non-

```

algorithm excess scaling;
begin
  preprocess;
   $\Delta := 2^{\lceil \log U \rceil}$ ;
  while  $\Delta \geq 1$  do
    begin ( $\Delta$ -scaling phase)
      while the network contains a node  $i$  with a large excess do
        begin
          among all nodes with a large excess, select a node  $i$  with
            the smallest distance label;
          perform  $\text{push/relabel}(i)$  while ensuring that no node excess exceeds  $\Delta$ ;
        end;
       $\Delta := \Delta/2$ ;
    end;
  end;

```

Figure 7.18 Excess scaling algorithm.

saturating, it sends $\min\{e(i), \Delta - e(j)\} \geq \Delta/2$ units of flow, proving the first part of the lemma. This push operation increases the excess of only node j . The new excess of node j is $e(j) + \min\{e(i), \Delta - e(j)\} \leq e(j) + \{\Delta - e(j)\} \leq \Delta$. So all the node excesses remain less than or equal to Δ . This proves the second part of the lemma. ♦

Lemma 7.21. *The excess scaling algorithm performs $O(n^2)$ nonsaturating pushes per scaling phase and $O(n^2 \log U)$ pushes in total.*

Proof. Consider the potential function $\Phi = \sum_{i \in N} e(i)d(i)/\Delta$. Using this potential function, we will establish the first assertion of the lemma. Since the algorithm performs $O(\log U)$ scaling phases, the second assertion is a consequence of the first. The initial value of Φ at the beginning of the Δ -scaling phase is bounded by $2n^2$ because $e(i)$ is bounded by Δ and $d(i)$ is bounded by $2n$. During the push/relabel(i) operation, one of the following two cases must apply:

Case 1. The algorithm is unable to find an admissible arc along which it can push flow. In this case the distance label of node i increases by $\epsilon \geq 1$ units. This relabeling operation increases Φ by at most ϵ units because $e(i) \leq \Delta$. Since for each i the total increase in $d(i)$ throughout the running of the algorithm is bounded by $2n$ (by Lemma 7.13), the total increase in Φ due to the relabeling of nodes is bounded by $2n^2$ in the Δ -scaling phase (actually, the increase in Φ due to node relabelings is at most $2n^2$ over *all* scaling phases).

Case 2. The algorithm is able to identify an arc on which it can push flow, so it performs either a saturating or a nonsaturating push. In either case, Φ decreases. A nonsaturating push on arc (i, j) sends at least $\Delta/2$ units of flow from node i to node j and since $d(j) = d(i) - 1$, after this operation decreases Φ by at least $\frac{1}{2}$ unit. Since the initial value of Φ at the beginning of a Δ -scaling phase is at most $2n^2$ and the increases in Φ during this scaling phase sum to at most $2n^2$ (from Case 1), the number of nonsaturating pushes is bounded by $8n^2$. ♦

This lemma implies a bound of $O(nm + n^2 \log U)$ on the excess scaling algorithm since we have already seen that all the other operations—such as saturating pushes, relabel operations, and finding admissible arcs—require $O(nm)$ time. Up to this point we have ignored the method needed to identify a node with the minimum distance label among nodes with excess more than $\Delta/2$. Making this identification is easy if we use a scheme similar to the one used in the highest-label preflow-push algorithm in Section 7.8 to find a node with the highest distance label. We maintain the lists $LIST(k) = \{i \in N : e(i) > \Delta/2 \text{ and } d(i) = k\}$, and a variable *level* that is a lower bound on the smallest index k for which $LIST(k)$ is nonempty. We identify the lowest-indexed nonempty list by starting at $LIST(\text{level})$ and sequentially scanning the higher-indexed lists. We leave as an exercise to show that the overall effort needed to scan the lists is bounded by the number of pushes performed by the algorithm plus $O(n \log U)$, so these computations are not a bottleneck operation. With this observation we can summarize our discussion as follows.

Theorem 7.22. *The excess scaling algorithm runs in $O(nm + n^2 \log U)$ time.* ♦

Algorithm	Running time	Features
Labeling algorithm	$O(nmU)$	<ol style="list-style-type: none"> 1. Maintains a feasible flow and augments flows along directed paths in the residual network from node s to node t. 2. Easy to implement and very flexible. 3. Running time is pseudopolynomial: the algorithm is not very efficient in practice.
Capacity scaling algorithm	$O(nm \log U)$	<ol style="list-style-type: none"> 1. A special implementation of the labeling algorithm. 2. Augments flows along paths from node s to node t with sufficiently large residual capacity. 3. Unlikely to be efficient in practice.
Successive shortest path algorithm	$O(n^2m)$	<ol style="list-style-type: none"> 1. Another special implementation of the labeling algorithm. 2. Augments flows along shortest directed paths from node s to node t in the residual network. 3. Uses distance labels to identify shortest paths from node s to node t. 4. Relatively easy to implement and very efficient in practice.
Generic preflow-push algorithm	$O(n^2m)$	<ol style="list-style-type: none"> 1. Maintains a pseudoflow; performs push/relabel operations at active nodes. 2. Very flexible; can examine active nodes in any order. 3. Relatively difficult to implement because an efficient implementation requires the use of several heuristics.
FIFO preflow-push algorithm	$O(n^3)$	<ol style="list-style-type: none"> 1. A special implementation of the generic preflow-push algorithm. 2. Examines active nodes in the FIFO order. 3. Very efficient in practice.
Highest-label preflow-push algorithm	$O(n^2\sqrt{m})$	<ol style="list-style-type: none"> 1. Another special implementation of the generic preflow-push algorithm. 2. Examines active nodes with the highest distance label. 3. Possibly the most efficient maximum flow algorithm in practice.
Excess scaling algorithm	$O(nm + n^2 \log U)$	<ol style="list-style-type: none"> 1. A special implementation of the generic preflow-push algorithm. 2. Performs push/relabel operations at nodes with sufficiently large excesses and, among these nodes, selects a node with the smallest distance label. 3. Achieves an excellent running time without using sophisticated data structures.

Figure 7.19 Summary of maximum flow algorithms.

7.10 SUMMARY

Building on the labeling algorithm described in Chapter 6, in this chapter we described several polynomial-time algorithms for the maximum flow problem. The labeling algorithm can perform as many as nU augmentations because each augmentation might carry a small amount of flow. We studied two natural strategies for reducing the number of augmentations and thus for improving the algorithm's running time; these strategies lead to the capacity scaling algorithm and the shortest augmenting path algorithm. One inherent drawback of these augmenting path algorithms is the computationally expensive operation of sending flows along paths. These algorithms might repeatedly augment flows along common path segments. The preflow-push algorithms that we described next overcome this drawback; we can conceive of them as sending flows along several paths simultaneously. In our development we considered both a generic implementation and several specific implementations of the preflow-push algorithm. The FIFO and highest-label preflow-push algorithms choose the nodes for pushing/relabeling in a specific order. The excess scaling algorithm ensures that the push operations, and subsequent augmentations, do not carry small amounts of flow (with "small" defined dynamically throughout the algorithm). Figure 7.19 summarizes the running times and basic features of these algorithms.

REFERENCE NOTES

The maximum flow problem is distinguished by the long succession of research contributions that have improved on the worst-case complexity of the best known algorithms. Indeed, no other network flow problem has witnessed as many incremental improvements. The following discussion provides a brief survey of selective improvements; Ahuja, Magnanti, and Orlin [1989, 1991] give a more complete survey of the developments in this field.

The labeling algorithm of Ford and Fulkerson [1956a] runs in pseudopolynomial time. Edmonds and Karp [1972] suggested two polynomial-time implementations of this algorithm. The first implementation, which augments flow along paths with the maximum residual capacity, performs $O(m \log U)$ iterations. The second implementation, which augments flow along shortest paths, performs $O(nm)$ iterations and runs in $O(nm^2)$ time. Independently, Dinic [1970] introduced a concept of shortest path networks (in number of arcs), called *layered networks*, and obtained an $O(n^2m)$ -time algorithm. Until this point all maximum flow algorithms were augmenting path algorithms. Karzanov [1974] introduced the first preflow-push algorithm on layered networks; he obtained an $O(n^3)$ algorithm. Shiloach and Vishkin [1982] described another $O(n^3)$ preflow-push algorithm for the maximum flow problem, which is a precursor of the FIFO preflow-push algorithm that we described in Section 7.7.

The capacity scaling described in Section 7.3 is due to Ahuja and Orlin [1991]; this algorithm is similar to the bit-scaling algorithm due to Gabow [1985] that we describe in Exercise 7.19. The shortest augmenting path algorithm described in Section 7.4 is also due to Ahuja and Orlin [1991]; this algorithm can be regarded as a variant of Dinic's [1970] algorithm and uses distance labels instead of layered networks.

Researchers obtained further improvements in the running times of the maximum flow algorithms by using distance labels instead of layered networks. Goldberg [1985] first introduced distance labels; by incorporating them in the algorithm of Shiloach and Vishkin [1982], he obtained the $O(n^3)$ -time FIFO implementation that we described in Section 7.7. The generic preflow-push algorithm and its highest-label preflow-push implementation that we described in Sections 7.8 are due to Goldberg and Tarjan [1986]. Using a dynamic tree data structure developed by Sleator and Tarjan [1983], Goldberg and Tarjan [1986] improved the running time of the FIFO implementation to $O(nm \log(n^2/m))$. Using a clever analysis, Cheriyan and Maheshwari [1989] show that the highest-label preflow-push algorithm in fact runs in $O(n^2 \sqrt{m})$ time. Our discussion in Section 7.7 presents a simplified proof of Cheriyan and Maheshwari approach. Ahuja and Orlin [1989] developed the excess scaling algorithm described in Section 7.9; this algorithm runs in $O(nm + n^2 \log U)$ time and obtains dramatic improvements over the FIFO and highest-label preflow-push algorithms without using sophisticated data structures. Ahuja, Orlin, and Tarjan [1989] further improved the excess scaling algorithm and obtained several algorithms: the best time bound of these algorithms is $O(nm \log(n \sqrt{\log U/m + 2}))$.

Cheriyan and Hagerup [1989] proposed a randomized algorithm for the maximum flow problem that has an expected running time of $O(nm)$ for all $m \geq n \log^2 n$. Alon [1990] developed a nonrandomized version of this algorithm and obtained a (deterministic) maximum flow algorithm that runs in (1) $O(nm)$ time for all $m = \Omega(n^{5/3} \log n)$, and (2) $O(nm \log n)$ time for all other values of n and m . Cheriyan, Hagerup, and Mehlhorn [1990] obtained an $O(n^3/\log n)$ algorithm for the maximum flow problem. Currently, the best available time bounds for solving the maximum flow problem are due to Alon [1990], Ahuja, Orlin, and Tarjan [1989], and Cheriyan, Hagerup, and Mehlhorn [1990].

Researchers have also investigated whether the worst-case bounds of the maximum flow algorithms are “tight” (i.e., whether algorithms achieve their worst-case bounds for some families of networks). Galil [1981] constructed a family of networks and showed that the algorithms of Edmonds and Karp [1972], Dinic [1970], Karzanov [1974], and a few other maximum flow algorithms achieve their worst-case bounds. Using this family of networks, it is possible to show that the shortest augmenting path algorithm also runs in $\Omega(n^2 m)$ time. Cheriyan and Maheshwari [1989] have shown that the generic preflow-push algorithm and its FIFO and the highest-label preflow-push implementations run in $\Omega(n^2 m)$, $\Omega(n^3)$, and $\Omega(n^2 \sqrt{m})$ times, respectively. Thus the worst-case time bounds of these algorithms are tight.

Several computational studies have assessed the empirical behavior of maximum flow algorithms. Among these, the studies by Imai [1983], Glover, Klingman, Mote, and Whitman [1984], Derigs and Meier [1989], and Ahuja, Kodialam, Mishra, and Orlin [1992] are noteworthy. These studies find that preflow-push algorithms are faster than augmenting path algorithms. Among the augmenting path algorithms, the shortest augmenting path algorithm is the fastest, and among the preflow-push algorithms, the performance of the highest-label preflow-push algorithm is the most attractive.

EXERCISES

- 7.1. Consider the network shown in Figure 7.20. The network depicts only those arcs with a positive capacity. Specify the residual network with respect to the current flow and compute exact distance labels in the residual network. Next change the arc flows (without changing the flow into the sink) so that the exact distance label of the source node (1) decreases by 1 unit; (2) increases by 1 unit.

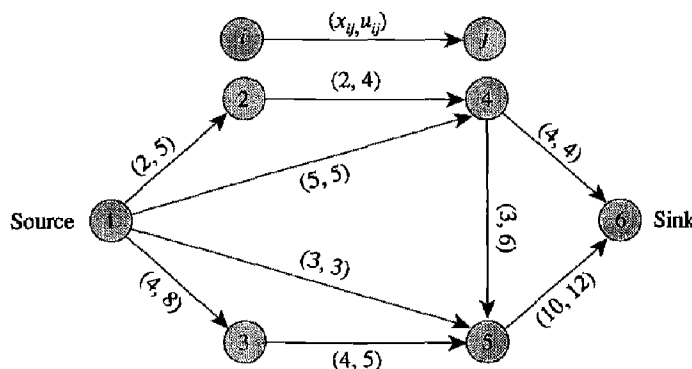


Figure 7.20 Example for Exercise 7.1.

- 7.2. Using the capacity scaling algorithm described in Section 7.3, find a maximum flow in the network given in Figure 7.21(b).

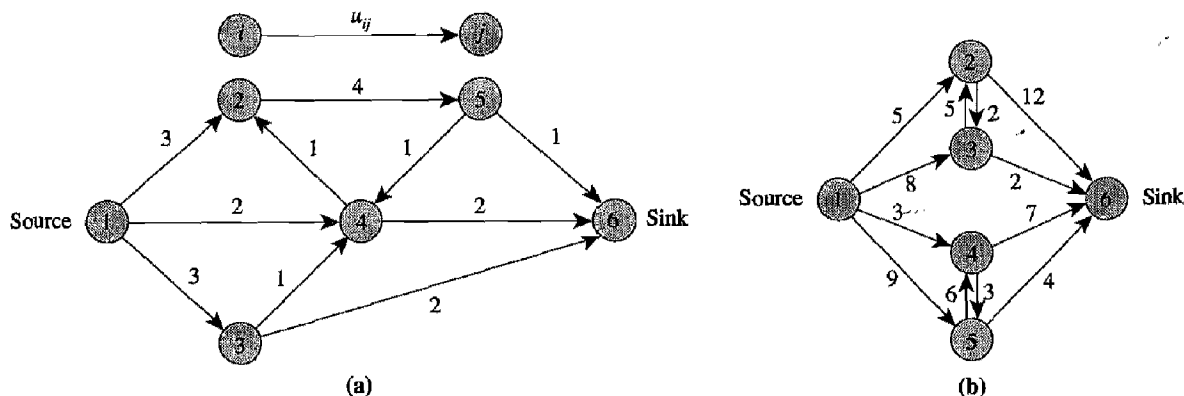


Figure 7.21 Examples for Exercises 7.2, 7.3, 7.5, and 7.6.

- 7.3. Using the shortest augmenting path algorithm, solve the maximum flow problem shown in Figure 7.21(a).
- 7.4. Solve the maximum flow problem shown in Figure 7.22 using the generic preflow-push algorithm. Incorporate the following rules to maintain uniformity of your computations: (1) Select an active node with the smallest index. [For example, if nodes 2 and 3 are active, select node 2.] (2) Examine the adjacency list of any node in the increasing order of the head node indices. [For example, if $A(1) = \{(1, 5), (1, 2), (1, 7)\}$, then examine arc (1, 2) first.] Show your computations on the residual networks encountered during the intermediate iterations of the algorithm.

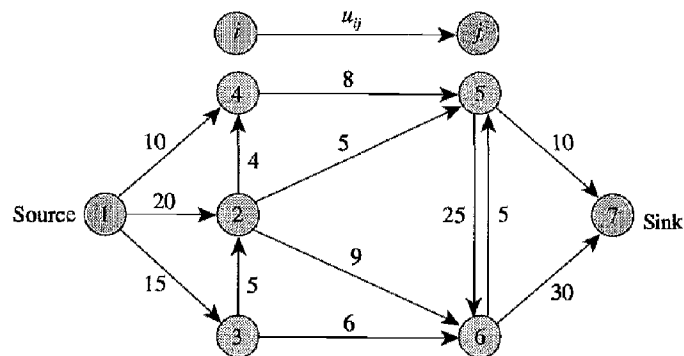


Figure 7.22 Example for Exercise 7.4.

- 7.5. Solve the maximum flow problem shown in Figure 7.21(a) using the FIFO preflow-push algorithm. Count the number of saturating and nonsaturating pushes and the number of relabel operations. Next, solve the same problem using the highest-label preflow-push algorithm. Compare the number of saturating pushes, nonsaturating pushes, and relabel operations with those of the FIFO preflow-push algorithm.
- 7.6. Using the excess scaling algorithm, determine a maximum flow in the network given in Figure 7.21(b).
- 7.7. **Most vital arcs.** We define a *most vital arc* of a network as an arc whose deletion causes the largest decrease in the maximum flow value. Either prove the following claims or show through counterexample that they are false.
- (a) A most vital arc is an arc with the maximum value of u_{ij} .
 - (b) A most vital arc is an arc with the maximum value of x_{ij} .
 - (c) A most vital arc is an arc with the maximum value of x_{ij} among arcs belonging to some minimum cut.
 - (d) An arc that does not belong to some minimum cut cannot be a most vital arc.
 - (e) A network might contain several most vital arcs.
- 7.8. **Least vital arcs.** A *least vital arc* in a network is an arc whose deletion causes the least decrease in the maximum flow value. Either prove the following claims or show that they are false.
- (a) Any arc with $x_{ij} = 0$ in any maximum flow is a least vital arc.
 - (b) A least vital arc is an arc with the minimum value of x_{ij} in a maximum flow.
 - (c) Any arc in a minimum cut cannot be a least vital arc.
- 7.9. Indicate which of the following claims are true or false. Justify your answer by giving a proof or by constructing a counterexample.
- (a) If the capacity of every arc in a network is a multiple of α , then in every maximum flow, each arc flow will be a multiple of α .
 - (b) In a network G , if the capacity of every arc increases by α units, the maximum flow value will increase by a multiple of α .
 - (c) Let v^* denote the maximum flow value of a given maximum flow problem. Let v' denote the flow into the sink node t at some stage of the preflow-push algorithm. Then $v^* - v' \leq \sum_{i \text{ is active}} e(i)$.
 - (d) By the flow decomposition theory, some sequence of at most $m + n$ augmentations would always convert any preflow into a maximum flow.
 - (e) In the excess scaling algorithm, $e_{\max} = \max\{e(i) : i \text{ is active}\}$ is a nonincreasing function of the number of push/relabel steps.
 - (f) The capacity of the augmenting paths generated by the maximum capacity augmenting path algorithm is nonincreasing.
 - (g) If each distance label $d(i)$ is a lower bound on the length of a shortest path from node i to node t in the residual network, the distance labels are valid.

- 7.10. Suppose that the capacity of every arc in a network is a multiple of α and is in the range $[0, \alpha K]$ for some integer K . Does this information improve the worst-case complexity of the labeling algorithm, FIFO preflow-push algorithm, and the excess scaling algorithm?
- 7.11. **Converting a maximum preflow to a maximum flow.** We define a *maximum preflow* x^0 as a preflow with the maximum possible flow into the sink.
- Show that for a given maximum preflow x^0 , some maximum flow x^* with the same flow value as x^0 , satisfies the condition that $x_{ij}^* \leq x_{ij}^0$ for all arcs $(i, j) \in A$. (*Hint:* Use flow decomposition.)
 - Suggest a labeling algorithm that converts a maximum preflow into a maximum flow in at most $n + m$ augmentations.
 - Suggest a variant of the shortest augmenting path algorithm that would convert a maximum preflow into a maximum flow in $O(nm)$ time. (*Hint:* Define distance labels from the source node and show that the algorithm will create at most m arc saturations.)
 - Suggest a variant of the highest-label preflow-push algorithm that would convert a maximum preflow into a maximum flow. Show that the running time of this algorithm is $O(nm)$. (*Hint:* Use the fact that we can delete an arc with zero flow from the network.)
- 7.12. (a) An arc is *upward critical* if increasing the capacity of this arc increases the maximum flow value. Does every network have an upward critical arc? Describe an algorithm for identifying all upward critical arcs in a network. The worst-case complexity of your algorithm should be substantially better than that of solving m maximum flow problems.
- (b) An arc is *downward critical* if decreasing the capacity of this arc decreases the maximum flow value. Is the set of upward critical arcs the same as the set of downward critical arcs? If not, describe an algorithm for identifying all downward critical arcs; analyze your algorithm's worst-case complexity.
- 7.13. Show that in the shortest augmenting path algorithm or in the preflow-push algorithm, if an arc (i, j) is inadmissible at some stage, it remains inadmissible until the algorithm relabels node i .
- 7.14. Apply the generic preflow-push algorithm to the maximum flow problem shown in Figure 7.23. Always examine a node with the smallest distance label and break ties in favor of a node with the smallest node number. How many saturating and nonsaturating pushes does the algorithm perform?

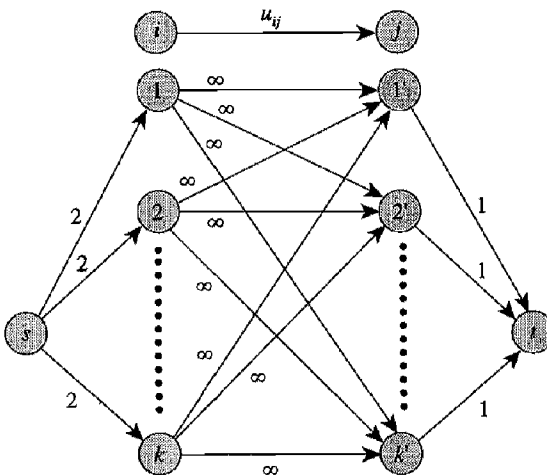


Figure 7.23 Example for Exercise 7.14.

- 7.15. Apply the FIFO preflow-push algorithm to the network shown in Figure 7.24. Determine the number of pushes as a function of the parameters W and L (correct within a constant factor). For a given value of n , what values of W and L produce the largest number of pushes?

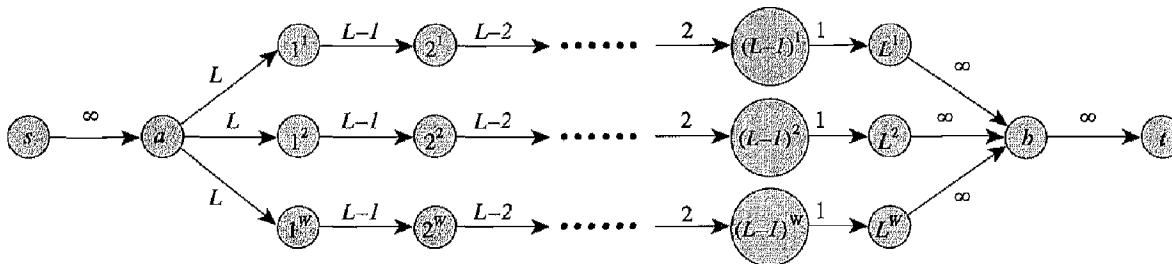


Figure 7.24 Example for Exercise 7.15.

- 7.16. Apply the highest-label preflow-push algorithm on the network shown in Figure 7.24. Determine the number of pushes as a function of the parameters W and L (correct within a constant factor). For a given n , what values of W and L produce the largest number of pushes?
- 7.17. Describe a more general version of the capacity scaling algorithm discussed in Section 7.3, one that scales Δ at each scaling phase by a factor of some integer number $\beta \geq 2$. Initially, $\Delta = \beta^{\lceil \log_\beta U \rceil}$ and each scaling phase reduces Δ by a factor of β . Analyze the worst-case complexity of this algorithm and determine the optimal value of β .
- 7.18. **Partially capacitated networks** (Ahuja and Orlin [1991]). Suppose that we wish to speed up the capacity scaling algorithm discussed in Exercise 7.17 for networks with some, but not all, arcs capacitated. Suppose that the network G has p arcs with finite capacities and $\beta = \max\{2, \lceil m/p \rceil\}$. Consider a version of the capacity scaling algorithm that scales Δ by a factor of β in each scaling phase. Show that the algorithm would perform at most $2m$ augmentations per scaling phase. [Hint: At the end of the Δ -scaling phase, the s - t cut in $G(\Delta)$ contains only arcs with finite capacities.] Conclude that the capacity scaling algorithm would solve the maximum flow problem in $O(m^2 \log_\beta U)$ time. Finally, show that this algorithm would run in $O(m^2)$ time if $U = O(n^k)$ for some k and $m = O(n^{1+\epsilon})$ for some $\epsilon > 0$.
- 7.19. **Bit-scaling algorithm** (Gabow [1985]). Let $K = \lceil \log U \rceil$. In the bit-scaling algorithm for the maximum flow problem works, we represent each arc capacity as a K -bit binary number, adding leading zeros if necessary to make each capacity K bits long. The problem P_k considers the capacity of each arc as the k leading bits. Let x_k^* denote a maximum flow and let v_k^* denote the maximum flow value in the problem P_k . The algorithm solves a sequence of problems $P_1, P_2, P_3, \dots, P_K$, using $2x_{k-1}^*$ as a starting solution for the problem P_k .
- (a) Show that $2x_{k-1}^*$ is feasible for P_k and that $v_k^* - 2v_{k-1}^* \leq m$.
- (b) Show that the shortest augmenting path algorithm for solving problem P_k , starting with $2x_{k-1}^*$ as the initial solution, requires $O(nm)$ time. Conclude that the bit-scaling algorithm solves the maximum flow problem in $O(nm \log U)$ time.
- 7.20. Using the potential function $\Phi = \max\{d(i) : i \text{ is active}\}$, show that the highest-label preflow-push algorithm performs $O(n^3)$ nonsaturating pushes.
- 7.21. The *wave algorithm*, which is a hybrid version of the highest-label and FIFO preflow-push algorithms, performs passes over active nodes. In each pass it examines *all* the active nodes in nonincreasing order of the distance labels. While examining a node, it pushes flow from a node until either its excess becomes zero or the node is relabeled. If during a pass, the algorithm relabels no node, it terminates; otherwise, in the next pass, it again examines active nodes in nonincreasing order of their new distance labels.

Discuss the similarities and differences between the wave algorithm with the highest label and FIFO preflow-push algorithms. Show that the wave algorithm runs in $O(n^3)$ time.

- 7.22. Several rules listed below are possible options for selecting an active node to perform the push/relabel operation in the preflow-push algorithm. Describe the data structure and the implementation details for each rule. Obtain the tightest possible bound on the numbers of pushes performed by the algorithm and the resulting running time of the algorithm.
- Select an active node with the minimum distance label.
 - Select an active node with the largest amount of excess.
 - Select an active node whose excess is at least 50 percent of the maximum excess at any node.
 - Select the active node that the algorithm had selected most recently.
 - Select the active node that the algorithm had selected least recently.
 - Select an active node randomly. (Assume that for any integer k , you can in $O(1)$ steps generate a random integer uniformly in the range $[1, k]$.)
- 7.23. In the excess scaling algorithm, suppose we require that each nonsaturating push pushes exactly $\Delta/2$ units of flow. Show how to modify the push/relabel step to meet this requirement. Does this modification change the worst-case running time of the algorithm?
- 7.24. In our development in this chapter we showed that the excess scaling algorithm performs $O(n^2 \log U^*)$ nonsaturating pushes if U^* is set equal to the largest arc capacity among the arcs emanating from the source node. However, we can often select smaller values of U^* and still show that the number of nonsaturating pushes is $O(n^2 \log U^*)$. Prove that we can also use the following values of U^* without affecting the worst-case complexity of the algorithm: (1) $U^* = \sum_{(s,j) \in A(s)} u_{sj} / |A(s)|$; (2) $U^* = v^{ub} / |A(s)|$ for any upper bound v^{ub} on the maximum flow value. (*Hint*: In the first scaling phase, set $\Delta = 2^{\lceil \log U^* \rceil}$ and forbid nodes except those adjacent to the source from having excess more than Δ .)
- 7.25. The excess scaling algorithm described in Section 7.9 scales excesses by a factor of 2. It starts with the value of the excess dominator Δ equal to the smallest power of 2 that is greater than or equal to U ; in every scaling phase, it reduces Δ by a factor of 2. An alternative is to scale the excesses by a factor of some integer number $\beta \geq 2$. This algorithm would run as follows. It would start with $\Delta = \beta^{\lceil \log_2 U \rceil}$; it would then reduce Δ by a factor of β in every scaling phase. In the Δ -scaling phase, we refer to a node with an excess of at least Δ/β as a node with a *large* excess. The algorithm pushes flow from a node with a large excess and among these nodes it chooses the node with the smallest distance label. The algorithm also ensures that no excess exceeds Δ . Determine the number of nonsaturating pushes performed by the algorithm as a function of n , β , and U . For what value of β would the algorithm perform the least number of nonsaturating pushes?
- 7.26. For any pair $[i, j] \in N \times N$, we define $\alpha[i, j]$ in the following manner: (1) if $(i, j) \in A$, then $\alpha[i, j]$ is the increase in the maximum flow value obtained by setting $u_{ij} = \infty$; and (2) if $(i, j) \notin A$, then $\alpha[i, j]$ is the increase in the maximum flow value obtained by introducing an infinite capacity arc (i, j) in the network.
- Show that $\alpha[i, j] \leq \alpha[s, j]$ and $\alpha[i, j] \leq \alpha[i, t]$.
 - Show that $\alpha[i, j] = \min\{\alpha[s, j], \alpha[i, t]\}$.
 - Show that we can compute $\alpha[i, j]$ for all node pairs by solving $O(n)$ maximum flow problems.
- 7.27. **Minimum cut with the fewest number of arcs.** Suppose that we wish to identify from among all minimum cuts, a minimum cut containing the least number of arcs. Show that if we replace u_{ij} by $u'_{ij} = mu_{ij} + 1$, the minimum cut with respect to the capacities u'_{ij} is a minimum cut with respect to the capacities u_{ij} containing the fewest number of arcs.

7.28. Parametric network feasibility problem. In a capacitated network G with arc capacities u_{ij} , suppose that the supply/demands of nodes are linear functions of time τ . Let each $b(i) = b^0(i) + \tau b^*(i)$ and suppose that $\sum_{i \in N} b^0(i) = 0$ and $\sum_{i \in N} b^*(i) = 0$. The network is currently (i.e., at time $\tau = 0$) able to fulfill the demands by the existing supplies but might not be able to do so at some point in future. You want to determine the largest integral value of τ up to which the network will admit a feasible flow. How would you solve this problem?

7.29. Source parametric maximum flow problem (Gallo, Grigoriadis, and Tarjan [1989]). In the *source parametric maximum flow problem*, the capacity of every source arc (s, j) is a nondecreasing linear function of a parameter λ (i.e., $u_{sj} = u_{sj}^0 + \lambda u_{sj}^*$ for some constant $u_{sj}^* \geq 0$); the capacity of every other arc is fixed, and we wish to determine a maximum flow for p values $0 = \lambda_1, \lambda_2, \dots, \lambda_p$ of the parameter λ . Assume that $\lambda_1 \leq \lambda_2 \leq \dots \leq \lambda_p$ and $p \leq n$. As an application of the source-parametric maximum flow problem, consider the following variation of Application 6.5. Suppose that processor 1 is a shared multiprogrammed system and processor 2 is a graphic processor dedicated to a single user. Suppose further that we can accurately determine the times required for processing modules on processor 2, but the times for processing modules on processor 1 are affected by a general work load on the processor. As the work load on processor 1 changes, the optimal distribution of modules between processor 1 and 2 changes. The source-parametric maximum flow problem determines these distributions for different work loads on processor 1.

Let $MF(\lambda)$ denote the maximum flow problem for a specific value of λ . Let $v(\lambda)$ denote the maximum flow value of $MF(\lambda)$ and let $[S(\lambda), \bar{S}(\lambda)]$ denote an associated minimum cut. Clearly, the zero flow is optimal for $MF(\lambda_1)$. Given an optimal flow $x(\lambda_k)$ of $MF(\lambda_k)$, we solve $MF(\lambda_{k+1})$ as follows: With $x(\lambda_k)$ as the starting flow and the corresponding distance labels as the initial distance labels, we perform a preprocess step by sending additional flow along the source arcs so that they all become saturated. Then we apply the FIFO preflow-push algorithm until the network contains no more active nodes. We repeat this process until we have solved $MF(\lambda_p)$.

- Show that the ending distance labels of $MF(\lambda_k)$ are valid distances for $MF(\lambda_{k+1})$ in the residual network $G(x(\lambda_k))$ after the preprocess step.
- Use the result in part (a) to show that overall [i.e., in solving all the problems $MF(\lambda_1), MF(\lambda_2), \dots, MF(\lambda_p)$], the algorithm performs $O(n^2)$ relabels, $O(nm)$ saturating pushes, and $O(n^3)$ nonsaturating pushes. Conclude that the FIFO preflow-push algorithm solves the source parametric maximum flow problem in $O(n^3)$ time, which is the same time required to solve a single maximum flow problem.
- Show that $v(\lambda_1) \leq v(\lambda_2) \leq \dots \leq v(\lambda_p)$ and some associated minimum cuts satisfy the nesting condition $S_1 \subseteq S_2 \subseteq \dots \subseteq S_p$.

7.30. Source-sink parametric maximum flow problem. In the source-sink parametric maximum flow problem, the capacity of every source arc is a *nondecreasing* linear function of a parameter λ and capacity of every sink arc is a *nonincreasing* linear function of λ , and we want to determine a maximum flow for several values of parameter $\lambda_1, \lambda_2, \dots, \lambda_p$, for $p \leq n$, that satisfy the condition $0 = \lambda_1 < \lambda_2 < \dots < \lambda_p$. Show how to solve this problem in a total of $O(n^3)$ time. (*Hint:* The algorithm is same as the one considered in Exercise 7.29 except that in the preprocess step if some sink arc has flow greater than its new capacity, we decrease the flow.)

7.31. Ryser's theorem. Let Q be a $p \times p$ matrix consisting of 0–1 elements. Let α denote the vector of row sums of Q and β denote the vector of column sums. Suppose that the rows and columns are ordered so that $\alpha_1 \geq \alpha_2 \geq \dots \geq \alpha_p$, and $\beta_1 \geq \beta_2 \geq \dots \geq \beta_p$.

- Show that the vectors α and β must satisfy the following conditions: (1) $\sum_{i=1}^p \alpha_i = \sum_{i=1}^p \beta_i$ and (2) $\sum_{i=1}^k \min(\alpha_i, k) \leq \sum_{i=1}^k \beta_i$, for all $k = 1, \dots, p$. [*Hint:* $\min(\alpha_i, k)$ is an upper bound on the sum of the first k components of row i .]

- (b) Given the nonnegative integer vector α and β , show how to formulate the problem of determining whether some 0–1 matrix Q has a row sum vector α and a column sum vector β as a maximum flow problem. Use the max-flow min-cut theorem to show that the conditions stated in part (a) are sufficient for the existence of such a matrix Q .

8

MAXIMUM FLOWS: ADDITIONAL TOPICS

*This was the most unkindest cut of all.
—Shakespeare in Julius Caesar Act III*

Chapter Outline

-
-
- 8.1 Introduction
 - 8.2 Flows in Unit Capacity Networks
 - 8.3 Flows in Bipartite Networks
 - 8.4 Flows in Planar Undirected Networks
 - 8.5 Dynamic Tree Implementations
 - 8.6 Network Connectivity
 - 8.7 All-Pairs Minimum Value Cut Problem
 - 8.8 Summary
-
-

8.1 INTRODUCTION

In all scientific disciplines, researchers are always making trade-offs between the generality and the specificity of their results. Network flows embodies these considerations. In studying minimum cost flow problems, we could consider optimization models with varying degrees of generality: for example, in increasing order of specialization, (1) general constrained optimization problems, (2) linear programs, (3) network flows, (4) particular network flow models (e.g., shortest path and maximum flow problems), and (5) the same models defined on problems with specific topologies and/or cost structures. The trade-offs in choosing where to study across the hierarchy of possible models is apparent. As models become broader, so does the range of their applications. As the models become more narrow, available results often become refined and more powerful. For example, as shown by our discussion in previous chapters, algorithms for shortest path and maximum flow problems have very attractive worst-case and empirical behavior. In particular, the computational complexity of these algorithms grows rather slowly in the number of underlying constraints (i.e., nodes) and decision variables (arcs). For more general linear programs, or even for more general minimum cost flow problems, the best algorithms are not nearly as good.

In considering what class of problems to study, we typically prefer models that are generic enough to be rich, both in applications and in theory. As evidenced by the coverage in this book, network flows is a topic that meets this criterion. Yet, through further specialization, we can develop a number of more refined results. Our study of shortest paths and maximum flow problems in the last four chapters

has illustrated this fact. Even within these more particular problem classes, we have seen the effect of further specialization, which has led to us to discover more efficient shortest path algorithms for models with nonnegative costs and for models defined on acyclic graphs. In this chapter we carry out a similar program for maximum flow problems. We consider maximum flow problems with both (1) specialized data, that is, networks with unit capacity arcs, and (2) specialized topologies, namely, bipartite and planar networks.

For general maximum flow problems, the labeling algorithm requires $O(nmU)$ computations and the shortest augmenting path algorithm requires $O(n^2m)$ computations. When applied to unit capacity networks, these algorithms are guaranteed to perform even better. Both require $O(nm)$ computations. We obtain this improvement simply because of the special nature of unit capacity networks. By designing specialized algorithms, however, we can improve even further on these results. Combining features of both the labeling algorithm and the shortest augmenting path algorithm, the unit capacity maximum flow algorithm that we consider in this chapter requires only $O(\min\{n^{2/3}m, m^{3/2}\})$ computations.

Network connectivity is an important application context for the unit capacity maximum flow problem. The arc connectivity between any two nodes of a network is the maximum number of arc-disjoint paths that connect these nodes; the arc connectivity of the network as a whole is the minimum arc connectivity between any pair of nodes. To determine this important reliability measure of a network, we could solve a unit capacity maximum flow problem between every pair of nodes, thus requiring $O(\min\{n^{2/3}m, m^{3/2}\})$ computations. As we will see in this chapter, by exploiting the special structure of the arc connectivity problem, we can reduce this complexity bound considerably—to $O(nm)$.

For networks with specialized bipartite and planar topologies, we can also obtain more efficient algorithms. Recall that bipartite networks are composed of two node sets, N_1 with n_1 nodes and N_2 with n_2 nodes. Assume that $n_1 \leq n_2$. For these problems we develop a specialization of the generic preflow-push algorithm that requires $O(n_1^2 m)$ instead of $O((n_1 + n_2)^2 m)$ time. Whenever the bipartite network is unbalanced in the sense that $n_1 \ll (n_1 + n_2) = n$, the new implementation has a much better complexity than the general preflow-push algorithm. Planar networks are those that we can draw on the plane so that no two arcs intersect each other. For this class of networks, we develop a specialized maximum flow algorithm that requires only $O(n \log n)$ computations.

In this chapter we also consider two other additional topics: a dynamic tree implementation and the all-pairs minimum value cut problem. Dynamic trees is a special type of data structure that permits us to implicitly send flow on paths of length n in $O(\log n)$ steps on average. By doing so we are able to reduce the computational requirement of the shortest augmenting path algorithm for maximum flows from $O(n^2m)$ to $O(nm \log n)$.

In some application contexts, we need to find the maximum flow between every pair of nodes in a network. The max-flow min-cut theorem shows that this problem is equivalent to finding the minimum cut separating all pairs of nodes. The most naive way to solve this problem would be to solve the maximum flow problem $n(n - 1)$ times, once between every pair of nodes. Can we do better? In Section 8.7 we show that how to exploit the relationship of the cut problems between various

node pairs to reduce the computational complexity of the all-pairs minimum value cut problem considerably in undirected networks. This algorithm requires solving only $(n - 1)$ minimum cut problems in undirected networks. Moreover, the techniques used in this development extend to a broader class of problems: they permit us to solve the all-pairs minimum cut problem for situations when the value of a cut might be different than the sum of the arc capacities across the cut.

The algorithms we examine in this chapter demonstrate the advantage of exploiting special structures to improve on the design of algorithms. This theme not only resurfaces on several other occasions in this book, but also is an important thread throughout the entire field of large-scale optimization. Indeed, we might view the field of large-scale optimization, and the field of network flows for that matter, as the study of theory and algorithms for exploiting special problem structure. In this sense this chapter is, in its orientation and overall approach, a microcosm of this entire book and of much of the field of optimization itself.

8.2 FLOWS IN UNIT CAPACITY NETWORKS

Certain combinatorial problems are naturally formulated as zero–one optimization models. When viewed as flow problems, these models yield networks whose arc capacities are all 1. We will refer to these networks as *unit capacity networks*. Frequently, it is possible to solve flow problems on these networks more efficiently than those defined on general networks. In this section we describe an efficient algorithm for solving the maximum flow problem on unit capacity networks. We subsequently refer to this algorithm as the *unit capacity maximum flow algorithm*.

In a unit capacity network, the maximum flow value is at most n , since the capacity of the s – t cut $\{s\}, S - \{s\}$ is at most n . The labeling algorithm therefore determines a maximum flow within n augmentations and requires $O(nm)$ effort. The shortest augmenting path algorithm also solves this problem in $O(nm)$ time since its bottleneck operation, which is the augmentation step, requires $O(nm)$ time instead of $O(n^2m)$ time. The unit capacity maximum flow algorithm that we describe is a hybrid version of these two algorithms. This unit capacity maximum flow algorithm is noteworthy because by combining features of both algorithms, it requires only $O(\min\{n^{2/3}m, m^{3/2}\})$ time, which is consistently better than the $O(nm)$ bound of either algorithm by itself.

The unit capacity maximum flow algorithm is a two-phase algorithm. In the first phase it applies the shortest augmenting path algorithm, although not until completion: rather, this phase terminates whenever the distance label of the source node satisfies the condition $d(s) \geq d^* = \min\{\lceil 2n^{2/3} \rceil, \lceil m^{1/2} \rceil\}$. Although the algorithm might terminate with a nonoptimal solution, the solution is probably nearly-optimal (its value is within d^* of the optimal flow value). In its second phase, the algorithm applies the labeling algorithm to convert this near-optimal flow into a maximum flow. As we will see, this two-phase approach works well for unit capacity networks because the shortest augmenting path algorithm obtains a near-optimal flow quickly (when augmenting paths are “short”) but then takes a long time to convert this solution into a maximum flow (when augmenting paths become “long”). It so happens that the labeling algorithm converts this near-optimal flow into a maximum flow far more quickly than the shortest augmenting path algorithm.

Let us examine the behavior of the shortest augmenting path algorithm for $d^* = \min\{\lceil 2n^{2/3} \rceil, \lceil m^{1/2} \rceil\}$. Suppose the algorithm terminates with a flow vector x' with a flow value equal to v' . What can we say about $v^* - v'$? (Recall that v^* denotes the maximum flow value.) We shall answer this question in two parts: (1) when $d^* = \lceil 2n^{2/3} \rceil$, and (2) when $d^* = \lceil m^{1/2} \rceil$.

Suppose that $d^* = \lceil 2n^{2/3} \rceil$. For each $k = 0, 1, 2, \dots, d^*$, let V_k denote the set of nodes with a distance label equal to k [i.e., $V_k = \{i \in N : d(i) = k\}$]. We refer to V_k as the set of nodes in the k th layer of the residual network. Consider the situation when each of the sets V_1, V_2, \dots, V_{d^*} is nonempty. It is possible to show that each arc (i, j) in the residual network $G(x')$ connects a node in the k th layer to a node in the $(k + 1)$ th layer for some k , for otherwise $d(i) > d(j) + 1$, which contradicts the distance label validity conditions (7.2). Therefore, for each $k = 1, 2, \dots, d^*$, the set of arcs joining the node sets V_k to V_{k-1} form an s - t cut in the residual network. In case one of the sets, say V_k , is empty, our discussion in Section 7.4 implies that the cut $[S, \bar{S}]$ defined by $S = V_{k+1} \cup V_{k+2} \cup \dots \cup V_{d^*}$ is a minimum cut.

Note that $|V_1| + |V_2| + \dots + |V_{d^*}| \leq n - 1$, because the sink node does not belong to any of these sets. We claim that the residual network contains at least two consecutive layers V_k and V_{k-1} , each with at most $n^{1/3}$ nodes. For if not, every alternate layer (say, V_1, V_3, V_5, \dots) must contain more than $n^{1/3}$ nodes and the total number of nodes in these layers would be strictly greater than $n^{1/3} d^*/2 \geq n$, leading to a contradiction. Consequently, $|V_k| \leq n^{1/3}$ and $|V_{k-1}| \leq n^{1/3}$ for some of the two layers V_k and V_{k-1} . The residual capacity of the s - t cut defined by the arcs connecting V_k to V_{k-1} is at most $|V_k| |V_{k-1}| \leq n^{2/3}$ (since at most one arc of unit residual capacity joins any pair of nodes). Therefore, by Property 6.2, $v^* - v' \leq n^{2/3} \leq d^*$.

Next consider the situation when $d^* = \lceil m^{1/2} \rceil$. The layers of nodes V_1, V_2, \dots, V_{d^*} define d^* s - t cuts in the residual network and these cuts are arc disjoint in $G(x)$. The sum of the residual capacities of these cuts is at most m since each arc contributes at most one to the residual capacity of any such cut. Thus some s - t cut must have residual capacity at most $\lceil m^{1/2} \rceil$. This conclusion proves that $v^* - v' \leq \lceil m^{1/2} \rceil = d^*$.

In both cases, whenever $d^* = \lceil 2n^{2/3} \rceil$ or $d^* = \lceil m^{1/2} \rceil$, we find that the first phase obtains a flow whose value differs from the maximum flow value by at most d^* units. The second phase converts this flow into a maximum flow in $O(d^*m)$ time since each augmentation requires $O(m)$ time and carries a unit flow. We now show that the first phase also requires $O(d^*m)$ time.

In the first phase, whenever the distance label of a node k exceeds d^* , this node never occurs as an intermediate node in any subsequent augmenting path since $d(k) < d(s) < d^*$. So the algorithm relabels any node at most d^* times. This observation gives a bound of $O(d^*n)$ on the number of retreat operations and a bound of $O(d^*m)$ on the time to perform the retreat operations. Consider next the augmentation time. Since each arc capacity is 1, flow augmentation over an arc immediately saturates that arc. During two consecutive saturations of any arc (i, j) , the distance labels of both the nodes i and j must increase by at least 2 units. Thus the algorithm can saturate any arc at most $\lfloor d^*/2 \rfloor$ times, giving an $O(d^*m)$ bound on the total time needed for flow augmentations. The total number of advance op-

erations is bounded by the augmentation time plus the number of retreat operations and is again $O(d^*m)$. We have established the following result.

Theorem 8.1. *The unit capacity maximum flow algorithm solves a maximum flow problem on unit capacity networks in $O(\min\{n^{2/3}m, m^{3/2}\})$ time.* ♦

The justification of this two-phase procedure should now be clear. If $d^* = \lceil 2n^{2/3} \rceil \leq m^{1/2}$, the preceding discussion shows that the shortest augmenting path algorithm requires $O(n^{2/3}m)$ computations to obtain a flow within $n^{2/3}$ of the optimal. If we allow this algorithm to run until it achieves optimality [i.e., until $d(s) \geq n$], the algorithm could require an additional $O((n - n^{2/3})m)$ time to convert this flow into an optimal flow. For $n = 1000$, these observations imply that if the algorithm achieves these bounds, it requires 10 percent of the time to send 90 percent of the maximum flow and the remaining 90 percent of the time to send 10 percent of the maximum flow. (Empirical investigations have observed a similar behavior in practice as well.) On the other hand, the use of labeling algorithm in the second phase establishes a maximum flow in $O(n^{2/3}m)$ time and substantially speeds up the overall performance of the algorithm.

Another special case of unit capacity networks, called *unit capacity simple networks*, also arises in practice and is of interest to researchers. For this class of unit capacity networks, every node in the network, except the source and sink nodes, has at most one incoming arc or at most one outgoing arc. The unit capacity maximum flow algorithm runs even faster for this class of networks. We achieve this improvement by setting $d^* = \lceil n^{1/2} \rceil$ in the algorithm.

Theorem 8.2. *The unit capacity maximum flow algorithm establishes a maximum flow in unit capacity simple networks in $O(n^{1/2}m)$ time.*

Proof. Consider the layers of nodes V_1, V_2, \dots, V_{d^*} at the end of the first phase. Note first that $d(s) > d^*$ since otherwise we could find yet another augmentation in Phase 1. Suppose that layer V_h contains the smallest number of nodes. Then $|V_h| \leq n^{1/2}$, since otherwise the number of nodes in all layers would be strictly greater than n . Let N' be the nodes in N with at most one outgoing arc. We define a cut $[S, N-S]$ as follows: $S = \{j: d(j) \geq h\} \cup \{j: d(j) = h \text{ and } j \in N'\}$. Since $d(s) > d^* \geq h$ and $d(t) = 0$, $[S, N-S]$ is an s - t cut. Each arc with residual capacity in the cut $[S, N-S]$ is either directed into a node in $V_h \cap (N-N')$ or else it is directed from a node in $V_h \cap N'$. Therefore, the residual capacity of the cut Q is at most $|V_h| \leq n^{1/2}$. Consequently, at the termination of the first phase, the flow value differs from the maximum flow value by at most $n^{1/2}$ units. Using arguments similar to those we have just used, we can now easily show that the algorithm would run in $O(n^{1/2}m)$ time. ♦

The proof of the Theorem 8.2 relies on the fact that only 1 unit of flow can pass through each node in the network (except the source and sink nodes). If we satisfy this condition but allow some arc capacities to be larger than 1, the unit capacity maximum flow algorithm would still require only $O(n^{1/2}m)$ time. Networks with this structure do arise on occasion; for example, we encountered this type of

network when we computed the maximum number of node-disjoint paths from the source node to the sink node in the proof of Theorem 6.8. Theorem 8.2 has another by-product: It permits us to solve the maximum bipartite matching problem in $O(n^{1/2}m)$ time since we can formulate this problem as a maximum flow problem on a unit capacity simple network. We study this transformation in Section 12.3.

8.3 FLOWS IN BIPARTITE NETWORKS

A *bipartite network* is a network $G = (N, A)$ with a node set N partitioned into two subsets N_1 and N_2 so that for every arc $(i, j) \in A$, either (1) $i \in N_1$ and $j \in N_2$, or (2) $i \in N_2$ and $j \in N_1$. We often represent a bipartite network using the notation $G = (N_1 \cup N_2, A)$. Let $n_1 = |N_1|$ and $n_2 = |N_2|$. Figure 8.1 gives an example of a bipartite network; in this case, we can let $N_1 = \{1, 2, 3, 9\}$ and $N_2 = \{4, 5, 6, 7, 8\}$.

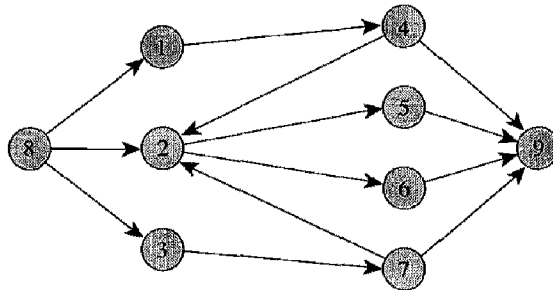


Figure 8.1 Bipartite network.

In this section we describe a specialization of the preflow-push algorithms that we considered in Chapter 7, but now adapt it to solve maximum flow problems on bipartite networks. The worst-case behavior of these special-purpose algorithms is similar to those of the original algorithms if the node sets N_1 and N_2 are of comparable size; the new algorithms are considerably faster than the original algorithms, however, whenever one of the sets N_1 or N_2 is substantially larger than the other. Without any loss of generality, we assume that $n_1 \leq n_2$. We also assume that the source node belongs to N_2 . [If the source node s belonged to N_1 , then we could create a new source node $s' \in N_2$, and we could add an arc (s', s) with capacity M for sufficiently large M .] As one example of the type of results we will obtain, we show that the specialization of the generic preflow-push algorithm solves a maximum flow problem on bipartite networks in $O(n_1^2 m)$ time. If $n_1 \ll (n_1 + n_2) = n$, the new implementation is considerably faster than the original algorithm.

In this section we examine only the generic preflow-push algorithm for bipartite networks; we refer to this algorithm as the *bipartite preflow-push algorithm*. The ideas we consider also apply in a straightforward manner to the FIFO, highest-label preflow-push and excess-scaling algorithms and yield algorithms with improved worst-case complexity. We consider these improvements in the exercises.

We first show that a slightly modified version of the generic preflow-push algorithm requires less than $O(n^2 m)$ time to solve problems defined on bipartite networks. To establish this result, we change the preprocess operation by setting $d(s) = 2n_1 + 1$ instead of $d(s) = n$. The modification stems from the observation

that any path in the residual network can have at most $2n_1$ arcs since every alternate node in the path must be in N_1 (because the residual network is also bipartite) and no path can repeat a node in N_1 . Therefore, if we set $d(s) = 2n_1 + 1$, the residual network will never contain a directed path from node s to node t , and the algorithm will terminate with a maximum flow.

Lemma 8.3. *For each node $i \in N$, $d(i) < 4n_1 + 1$.*

Proof. The proof is similar to that of Lemma 7.12. ◆

The following result is a direct consequence of this lemma.

Lemma 8.4.

- (a) *Each distance label increases at most $O(n_1)$ times. Consequently, the total number of relabel operations is $O(n_1(n_1 + n_2)) = O(n_1n_2)$.*
- (b) *The number of saturating pushes is $O(n_1m)$.*

Proof. The proofs are similar to those of Lemmas 7.13 and 7.14. ◆

It is possible to show that the results of Lemma 8.4 yield a bound of $O((n_1m)n)$ on the number of nonsaturating pushes, as well as on the complexity of the generic preflow-push algorithm. Instead of considering the details of this approach, we next develop a modification of the generic algorithm that runs in $O(n_1^2m)$ time.

This modification builds on the following idea. To bound the nonsaturating pushes of any preflow-push algorithm, we typically use a potential function defined in terms of all the active nodes in the network. If every node in the network can be active, the algorithm will perform a certain number of nonsaturating pushes. However, if we permit only the nodes in N_1 to be active, because $n_1 \leq n$ we can obtain a tighter bound on the number of nonsaturating pushes. Fortunately, the special structure of bipartite networks permits us to devise an algorithm that always manages to keep the nodes in N_2 inactive. We accomplish this objective by starting with a solution whose only active nodes are in N_1 , and by performing pushes of length 2; that is, we push flow over two consecutive admissible arcs so that any excess always returns to a node in N_1 and no node in N_2 ever becomes active.

Consider the residual network of a bipartite network given in Figure 8.2(a), with node excesses and distance labels displayed next to the nodes, and residual capacities displayed next to the arcs. The bipartite preflow-push algorithm first pushes flow from node 1 because it is the only active node in the network. The algorithm then identifies an admissible arc emanating from node 1. Suppose that it selects the arc (1, 3). Since we want to find a path of length 2, we now look for an admissible arc emanating from node 3. The arc (3, 2) is one such arc. We perform a push on this path, pushing $\delta = \min\{e(1), r_{13}, r_{32}\} = \min\{6, 5, 4\} = 4$ units. This push saturates arc (3, 2), and completes one iteration of the algorithm. Figure 8.2(b) gives the solution at this point.

In the second iteration, suppose that the algorithm again selects node 1 as an active node and arc (1, 3) as an admissible arc emanating from this node. We would also like to find an admissible arc emanating from node 3, but the network has none. So we relabel node 3. As a result of this relabel operation, arc (1, 3) becomes in-

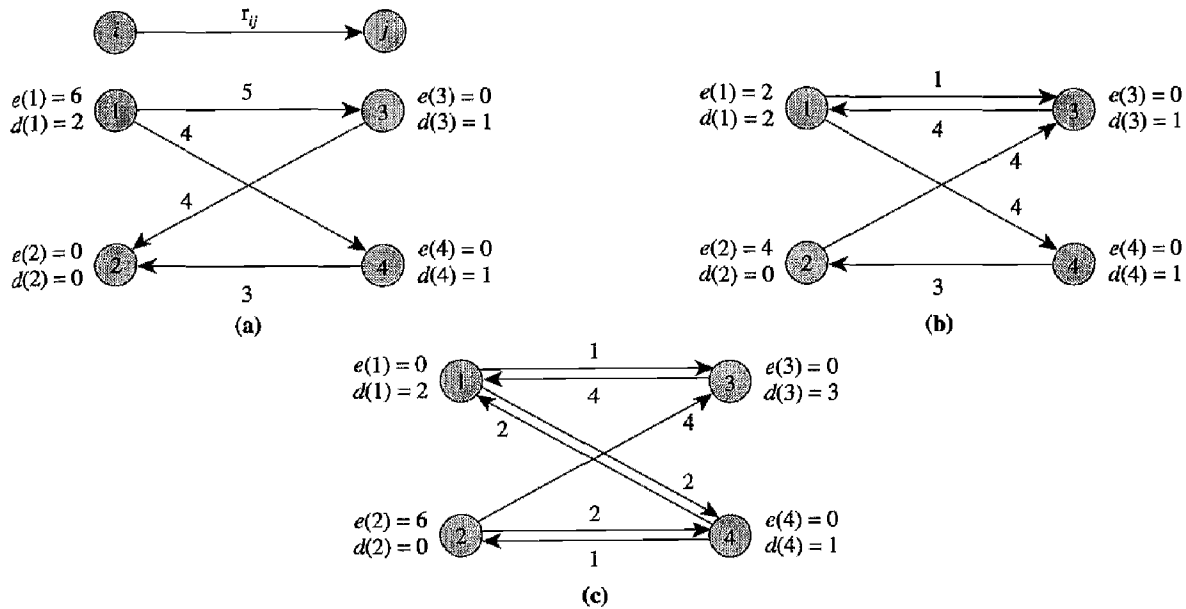


Figure 8.2 Illustrating the bipartite preflow-flow algorithm.

admissible. This operation completes the second iteration and Figure 8.2(b) gives the solution at this stage except that $d(3)$ is 3 instead of 1.

Suppose that the algorithm again selects node 1 as an active node in the third iteration. Then it selects the two consecutive admissible arcs $(1, 4)$ and $(4, 2)$, and pushes $\delta = \min\{e(1), r_{14}, r_{42}\} = \min\{2, 4, 3\} = 2$ units of flow over these arcs. This push is nonsaturating and eliminates the excess at node 1. Figure 8.2(c) depicts the solution at this stage.

As we have illustrated in this numerical example, the bipartite preflow-push algorithm is a simple generalization of the generic preflow-push algorithm. The bipartite algorithm is the same as the generic algorithm given in Figure 7.12 except that we replace the procedure `push/relabel(i)` by the procedure given in Figure 8.3.

```

procedure bipartite push/relabel( $i$ );
begin
  if the residual network contains an admissible arc  $(i, j)$  then
    if the residual network contains an admissible arc  $(j, k)$  then
      push  $\delta = \min\{e(i), r_{ij}, r_{jk}\}$  units of flow over the path  $i \rightarrow j \rightarrow k$ 
    else replace  $d(j)$  by  $\min\{d(k) + 1 : (j, k) \in A(j) \text{ and } r_{jk} > 0\}$ 
    else replace  $d(i)$  by  $\min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\}$ ;
end;

```

Figure 8.3 Push/relabel operation for bipartite networks.

Lemma 8.5. *The bipartite preflow-push algorithm performs $O(n_1^2 m)$ non-saturating pushes and runs in $O(n_1^2 m)$ time.*

Proof. The proof is same as that of Lemma 7.15. We consider the potential function $\Phi = \sum_{i \in I} d(i)$ whose index set I is the set of active nodes. Since we allow only the nodes in N_1 to be active, and $d(i) \leq 4n_1$ for all $i \in N_1$, the initial value of

Φ is at most $4n_1^2$. Let us observe the effect of executing the procedure `bipartite push/relabel(i)` on the potential function Φ . The procedure produces one of the following four outcomes: (1) it increases the distance label of node i ; (2) it increases the distance label of a node $j \in N_2$; (3) it pushes flow over the arcs (i, j) and (j, k) , saturating one of these two arcs; or (4) it performs a nonsaturating push. In case 1, the potential function Φ increases, but the total increase over all such iterations is only $O(n_1^2)$. In case 2, Φ remains unchanged. In case 3, Φ can increase by as much as $4n_1 + 1$ units since a new node might become active; Lemma 8.4 shows that the total increase over all iterations is $O(n_1^2 m)$. Finally, a nonsaturating push decreases the potential function by at least 2 units since it makes node i inactive, can make node k newly active and $d(k) = d(i) - 2$. This fact, in view of the preceding arguments, implies that the algorithm performs $O(n_1^2 m)$ nonsaturating pushes. Since all the other operations, such as the relabel operations and finding admissible arcs, require only $O(n_1 m)$ time, we have established the theorem. \blacklozenge

We complete this section by giving two applications of the maximum flow problem on bipartite networks with $n_1 \leq n_2$.

Application 8.1 Baseball Elimination Problem

At a particular point in the baseball season, each of $n + 1$ teams in the American League, which we number as $0, 1, \dots, n$, has played several games. Suppose that team i has won w_i of the games that it has already played and that g_{ij} is the number of games that teams i and j have yet to play with each other. No game ends in a tie. An avid and optimistic fan of one of the teams, the Boston Red Sox, wishes to know if his team still has a chance to win the league title. We say that we can *eliminate* a specific team 0, the Red Sox, if for every possible outcome of the unplayed games, at least one team will have more wins than the Red Sox. Let w_{\max} denote w_0 plus the total number of games team 0 has yet to play, which, in the best of all possible worlds, is the number of victories the Red Sox can achieve. Then we cannot eliminate team 0 if in some outcome of the remaining games to be played throughout the league, w_{\max} is at least as large as the possible victories of every other team. We want to determine whether we can or cannot eliminate team 0.

We can transform this baseball elimination problem into a feasible flow problem on a bipartite network with two sets with n_1 and $n_2 = \Omega(n_1^2)$. As discussed in Section 6.2, we can represent the feasible flow problem as a maximum flow problem, as shown in Figure 8.4. The maximum flow network associated with this problem contains n team nodes 1 through n , $n(n - 1)/2$ game nodes of the type $i-j$ for each $1 \leq i \leq j \leq n$, and source node s . Each game node $i-j$ has two incoming arcs $(i, i-j)$ and $(j, i-j)$, and the flows on these arcs represent the number of victories for team i and team j , respectively, among the additional g_{ij} games that these two teams have yet to play against each other (which is the required flow into the game node $i-j$). The flow x_{si} on the source arc (s, i) represents the total number of additional games that team i wins. We cannot eliminate team 0 if this network contains a feasible flow x satisfying the conditions

$$w_{\max} \geq w_i + x_{si} \quad \text{for all } i = 1, \dots, n,$$

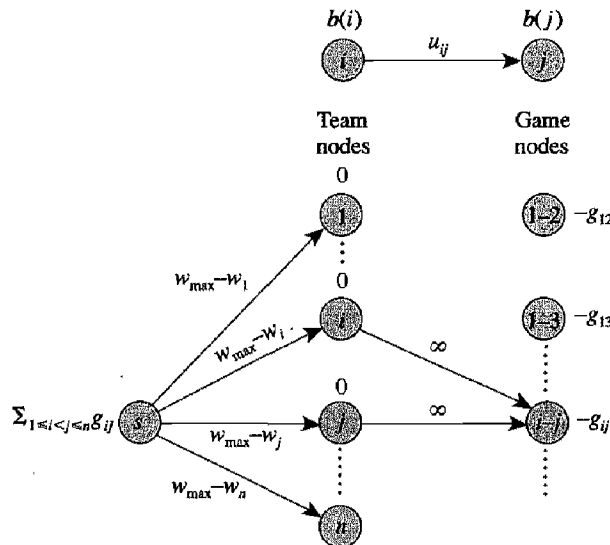


Figure 8.4 Network formulation of the baseball elimination problem.

which we can rewrite as

$$x_{si} \leq w_{\max} - w_i \quad \text{for all } i = 1, \dots, n.$$

This observation explains the capacities of arcs shown in the figure. We have thus shown that if the feasible flow problem shown in Figure 8.4 admits a feasible flow, we cannot eliminate team 0; otherwise, we can eliminate this team and our avid fan can turn his attention to other matters.

Application 8.2 Network Reliability Testing

In many application contexts, we need to test or monitor the arcs of a network (e.g., the tracks in a rail network) to ensure that the arcs are in good working condition. As a practical illustration, suppose that we wish to test each arc $(i, j) \in A$ in an undirected communication network $G = (N, A)$ α_{ij} times; due to resource limitations, however, each day we can test at most β_j arcs incident to any communication node $j \in N$. The problem is to find a schedule that completes the testing of all the arcs in the fewest number of days.

We solve this problem on a bipartite network $G' = (\{s\} \cup \{t\} \cup N_1 \cup N_2, A')$ defined as follows: The network contains a node $i \in N_1$ for every node $i \in N$ in the communication network and a node $i-j \in N_2$ for every arc $(i, j) \in A$ in the communication network. Each $i-j$ node has two incoming arcs from the nodes in N_1 , one from node i and the other from node j ; all these arcs have infinite capacity. The source node s is connected to every node $i \in N_1$ with an arc of capacity $\lambda\beta_j$, and every node $i-j \in N_2$ is connected to the sink node t with an arc of capacity α_{ij} . The reliability testing problem is to determine the smallest integral value of the days λ so that the maximum flow in the network saturates all the sink arcs. We can solve this problem by performing binary search on λ and solving a maximum flow problem at each search point. In these maximum flow problem, $|N_1| = n$ and $|N_2| = m$, and m can be as large as $n(n-1)/2$.

8.4 FLOWS IN PLANAR UNDIRECTED NETWORKS

A network is said to be *planar* if we can draw it in a two-dimensional (Euclidean) plane so that no two arcs cross (or intersect each other); that is, we allow the arcs to touch one another only at the nodes. Planar networks are an important special class of networks that arise in several application contexts. Because of the special structure of planar networks, network flow algorithms often run faster on these networks than they do on more general networks. Indeed, several network optimization problems are NP-complete on general networks (e.g., the maximum cut problem) but can be solved in polynomial time on planar networks. In this section we study some properties of planar networks and describe an algorithm that solves a maximum flow problem in planar networks in $O(n \log n)$ time. In this section we restrict our attention to undirected networks. We remind the reader that the undirected networks we consider contain at most one arc between any pair i and j of nodes. The capacity u_{ij} of arc (i, j) denotes the maximum amount that can flow from node i to node j or from node j to node i .

Figure 8.5 gives some examples of planar networks. The network shown in Figure 8.5(a) does not appear to be planar because arcs $(1, 3)$ and $(2, 4)$ cross one

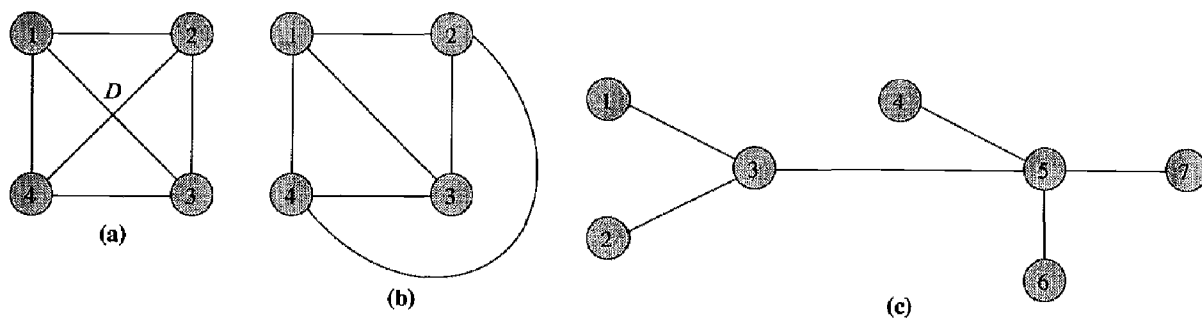


Figure 8.5 Instances of planar networks.

another at the point D , which is not a node. But, in fact, the network is planar because, as shown in Figure 8.5(b), we can redraw it, maintaining the network structure (i.e., node, arc structure), so that the arcs do not cross. For some networks, however, no matter how we draw them, some arcs will always cross. We refer to such networks as *nonplanar*. Figure 8.6 gives two instances of nonplanar networks. In both instances we could draw all but one arc without any arcs intersecting; if we add the last arc, though, at least one intersection is essential. Needless to say, determining whether a network is planar or not a straightforward task. However,

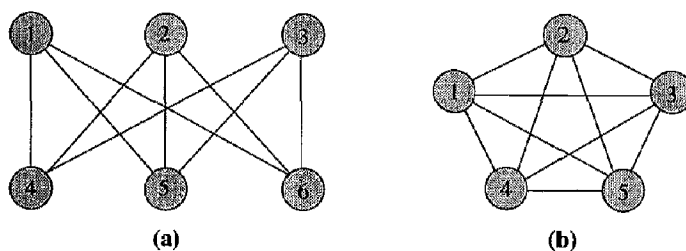


Figure 8.6 Instances of two nonplanar graphs.

researchers have developed very efficient algorithms (in fact, linear time algorithms) for testing the planarity of a network. Several theoretical characterizations of planar networks are also available.

Let $G = (N, A)$ be a planar network. A *face* z of G is a region of the (two-dimensional) plane bounded by arcs that satisfies the condition that any two points in the region can be connected by a continuous curve that meets no nodes and arcs. It is possible to draw a planar graph in several ways and each such representation might have a different set of faces. The *boundary* of a face z is the set of all arcs that enclose it. It is convenient to represent the boundary of a face by a cycle. Observe that each arc in the network belongs to the boundary of at most two faces. Faces z and z' are said to be *adjacent* if their boundaries contain a common arc. If two faces touch each other only at a node, we do not consider them to be adjacent. The network shown in Figure 8.5(b) illustrates these definitions. This network has four faces. The boundaries 1–2–3–1, 1–3–4–1, and 2–4–3–2 define the first three faces of the network. The fourth face is unbounded and consists of the remaining region; its boundary is 1–2–4–1. In Figure 8.5(b) each face is adjacent to every other face. The network shown in Figure 8.5(c) is a very special type of planar network. It has one unbounded face and its boundary includes all the arcs.

Next we discuss two well-known properties of planar networks.

Property 8.6 (Euler's Formula). *If a connected planar network has n nodes, m arcs, and f faces, then $f = m - n + 2$.*

Proof. We prove this property by performing induction on the value of f . For $f = 1$, $m = n - 1$, because a connected graph with just one face (which is the unbounded face) must be a spanning tree. Now assume, inductively, that Euler's formula is valid for every graph with k or fewer faces; we prove that the formula is valid for every graph with $k + 1$ faces. Consider a graph G with $k + 1$ faces and n nodes. We select any arc (i, j) that belongs to two faces, say z_1 and z_2 (show that the network always contains such an arc!). If we delete this arc from G , the two faces z_1 and z_2 merge into a single face. The resulting graph G' has m arcs, k faces, and n nodes, and by the induction hypothesis, $k = m - n + 2$. Therefore, if we reintroduce the arc (i, j) into G' , we see that $k + 1 = (m + 1) - n + 2$, so Euler's formula remains valid. We have thus completed the inductive step and established Euler's formula in general. ♦

Property 8.7. *In a planar network, $m < 3n$.*

Proof. We prove this property by contradiction. Suppose that $m \geq 3n$. Alternatively,

$$n \geq m/3. \quad (8.1)$$

We next obtain a relationship between f and m . Since the network contains no parallel arcs, the boundary of each face contains at least three arcs. Therefore, if we traverse the boundaries of all the faces one by one, we traverse at least $3f$ arcs. Now notice that we would have traversed each arc in the network at most twice because it belongs to the boundaries of at most two faces. These observations

show that $3f \leq 2m$. Alternatively,

$$f \leq 2m/3. \quad (8.2)$$

Using (8.1) and (8.2) in the formula $f = m - n + 2$, we obtain

$$2 = n - m + f \leq m/3 - m + 2m/3 = 0, \quad (8.3)$$

which is a contradiction. ♦

Property 8.7 shows that every planar graph is very sparse [i.e., $m = O(n)$]. This result, by itself, improves the running times for most network flow algorithms. For instance, as shown in Section 8.5, the shortest augmenting path algorithm for the maximum flow problem, implemented using the dynamic tree data structure, runs in $O(nm \log n)$ time. For planar networks, this time bound becomes $O(n^2 \log n)$. We can, in fact, develop even better algorithms by using the special properties of planar networks. To illustrate this point, we prove some results that apply to planar networks, but not to nonplanar networks. We show that we can obtain a minimum cut and a maximum flow for any planar network in $O(n \log n)$ time by solving a shortest path problem.

Finding Minimum Cuts Using Shortest Paths

Planar networks have many special properties. In particular, every connected planar network $G = (N, A)$ has an associated “twin” planar network $G^* = (N^*, A^*)$, which we refer to as the *dual* of G . We construct the dual G^* for a given graph G as follows. We first place a node f^* inside each face f of G . Each arc in G has a corresponding arc in G^* . Every arc (i, j) in G belongs to the boundaries of either (1) two faces, say f_1 and f_2 ; or (2) one face, say f_1 . In case 1, G^* contains the arc (f_1^*, f_2^*) ; in case 2, G^* contains the loop (f_1^*, f_1^*) . Figure 8.7, which illustrates this construction, depicts the dual network by dashed lines.

For notational convenience, we refer to the original network G as the *primal network*. The number of nodes in the dual network equals the number of faces in the primal network, and conversely, the number of faces in the dual network equals the number of nodes in the primal network. Both the primal and dual networks have

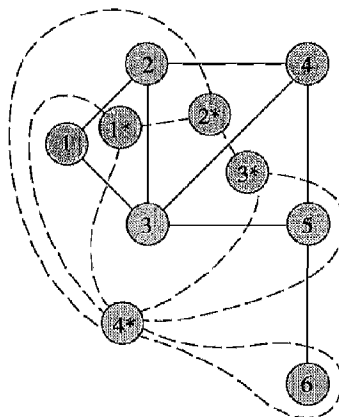


Figure 8.7 Constructing the dual of a planar network.

the same number of arcs. Furthermore, the dual of the dual network is the primal network. It is easy to show that a cycle in the dual network defines a cut in the primal network, and vice versa. For example, the cycle $4^*-1^*-2^*-4^*$ in the dual network shown in Figure 8.7 [with $(4^*, 1^*)$ denoting the arc from 4^* to 1^* that also passes through arc $(1, 2)$], defines the cut $\{(2, 1), (2, 3), (2, 4)\}$ in the primal network.

Our subsequent discussion in this section applies to a special class of planar networks known as s - t planar networks. A planar network with a source node s and a sink node t is called s - t planar if nodes s and t both lie on the boundary of the unbounded face. For example, the network shown in Figure 8.8(a) is s - t planar if $s = 1$ and $t = 8$; however, it is not s - t planar if (1) $s = 1$ and $t = 6$, or (2) $s = 3$ and $t = 8$.

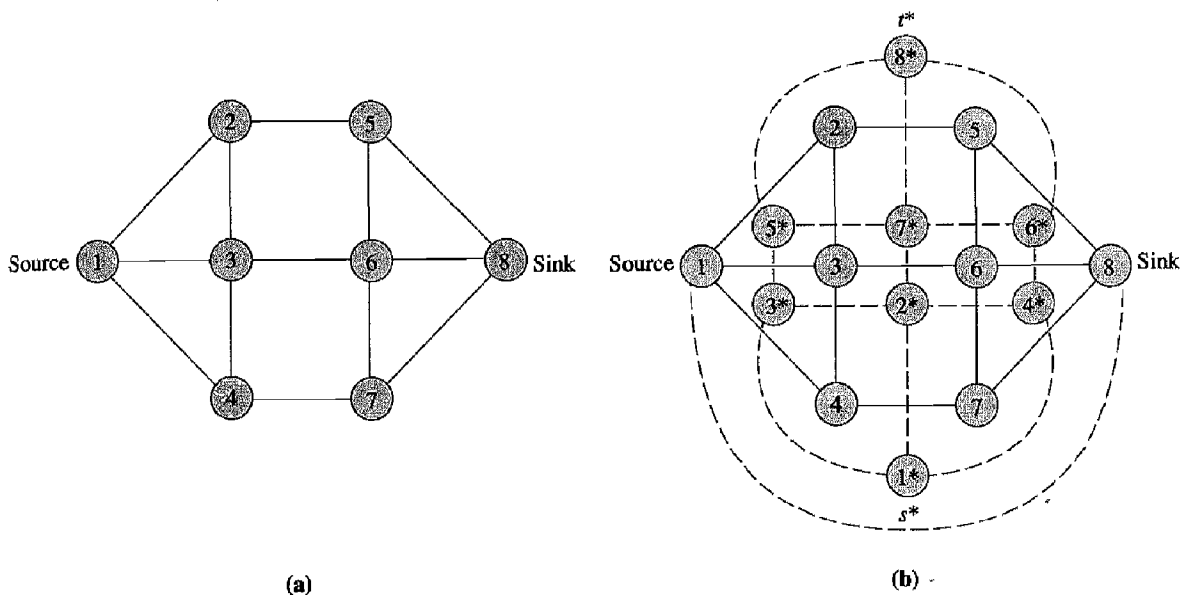


Figure 8.8 Establishing a relationship between cuts and paths: (a) s - t planar network; (b) corresponding dual network.

We now show how to transform a minimum cut problem on an s - t planar network into a shortest path problem. In the given s - t planar network, we first draw a new arc joining the nodes s and t so that the arc stays within the unbounded face of the network [see Figure 8.8(b)]; this construction creates a new face of the network, which we call the *additional face*, but maintains the network's planarity. We then construct the dual of this network; we designate the node corresponding to the additional face as the dual source s^* and the node corresponding to the unbounded face as the dual sink t^* . We set the cost of an arc in the dual network equal to the capacity of the corresponding arc in the primal network. The dual network contains the arc (s^*, t^*) which we delete from the network. Figure 8.8(b) shows this construction: the dashed lines are the arcs in the dual network. It is easy to establish a one-to-one correspondence between s - t cuts in the primal network and paths from node s^* to node t^* in the dual network; moreover, the capacity of the cut equals

the cost of the corresponding path. Consequently, we can obtain a minimum s - t cut in the primal network by determining a shortest path from node s^* to node t^* in the dual network.

In the preceding discussion we showed that by solving a shortest path problem in the dual network, we can identify a minimum cut in a primal s - t planar network. Since we can solve the shortest path problem in the dual network in $O(m \log n) = O(n \log n)$ using the binary heap implementation of Dijkstra's algorithm (see Section 4.7), this development provides us with an $O(n \log n)$ algorithm for identifying a minimum cut in a planar network. Notice that this bound is substantially better than the one we would obtain for a general network. We now give a generalization of this result, obtaining a rather surprising result that the shortest path distances in the dual network provide a maximum flow in the primal network.

Let $d(j^*)$ denote the shortest path distance from node s^* to node j^* in the dual network. Recall from Section 5.2 that the shortest path distances satisfy the following conditions:

$$d(j^*) \leq d(i^*) + c_{i^*j^*} \quad \text{for each } (i^*, j^*) \in A^*. \quad (8.4)$$

Each arc (i, j) in the primal network corresponds to an arc (i^*, j^*) in the dual network. Let us define a function x_{ij} for each $(i, j) \in A$ in the following manner:

$$x_{ij} = d(j^*) - d(i^*). \quad (8.5)$$

Note that $x_{ij} = -x_{ji}$. Now notice that the network G is undirected so that the arc set A contains both the arc (i, j) and the arc (j, i) . Hence we can regard a negative flow on arc (j, i) as a positive flow on arc (i, j) . Consequently, the flow vector x will always nonnegative.

The expressions (8.5) and (8.4) imply that

$$x_{ij} = d(j^*) - d(i^*) \leq c_{i^*j^*}. \quad (8.6)$$

Therefore, the flow x satisfies the arc capacity constraints. We next show that x also satisfies the mass balance constraints. Each node k in G , except node s and node t , defines a cut $Q = [\{k\}, N - \{k\}]$ consisting of all of the arcs incident to that node. The arcs in G^* corresponding to arcs in Q define a cycle, say W^* . For example, in Figure 8.7, the cycle corresponding to the cut for $k = 3$ is $1^*-2^*-3^*-4^*-1^*$. Clearly,

$$\sum_{(i^*, j^*) \in W^*} (d(j^*) - d(i^*)) = 0, \quad (8.7)$$

because the terms cancel each other. Using (8.5) in (8.7) shows that

$$\sum_{(i, j) \in Q} x_{ij} = 0,$$

which implies that inflow equals outflow at node k . Finally, we show that the flow x is a maximum flow. Let P^* be a shortest path from node s^* to node t^* in G^* . The definition of P^* implies that

$$d(j^*) - d(i^*) = c_{i^*j^*} \quad \text{for each } (i^*, j^*) \in P^*. \quad (8.8)$$

The arcs corresponding to P^* define an s - t cut Q in the primal network. Using (8.5) in expression (8.8) and using the fact that $c_{ij}^* = u_{ij}$, we get

$$x_{ij} = u_{ij} \quad \text{for each } (i, j) \in Q. \quad (8.9)$$

Consequently, the flow saturates all the arcs in an s - t cut and must be a maximum flow. The following theorem summarizes our discussion.

Theorem 8.8. *It is possible to determine a maximum flow in an s - t planar network in $O(n \log n)$ time.* ♦

8.5 DYNAMIC TREE IMPLEMENTATIONS

A dynamic tree is an important data structure that researchers have used extensively to improve the worst-case complexity of several network algorithms. In this section we describe the use of this data structure for the shortest augmenting path algorithm. We do not describe how to actually implement the dynamic tree data structure; rather, we show how to use this data structure as a “black box” to improve the computational complexity of certain algorithms. Our objective is to familiarize readers with this important data structure and enable them to use it as a black box module in the design of network algorithms.

The following observation serves as a motivation for the dynamic tree structure. The shortest augmenting path algorithm repeatedly identifies a path consisting solely of admissible arcs and augments flows on these paths. Each augmentation saturates some arcs on this path, and by deleting all the saturated arcs from this path we obtain a set of *path fragments*: sets of partial paths of admissible arcs. The path fragments contain valuable information. If we reach a node in any of these path fragments using any augmenting path, we know that we can immediately extend the augmenting path along the path fragment. The standard implementation of the shortest augmenting path algorithm discards this information and possibly regenerates it again at future steps. The dynamic tree data structure cleverly stores these path fragments and uses them later to identify augmenting paths quickly.

The dynamic tree data structure maintains a collection of node-disjoint rooted trees, each arc with an associated value, called *val*. See Figure 8.9(a) for an example of the node-disjoint rooted trees. Each rooted tree is a directed in-tree with a unique root. We refer to the nodes of the tree by using the terminology of a predecessor-successor (or parent-child) relationship. For example, node 5 is the predecessor (parent) of nodes 2 and 3, and nodes 9, 10, and 11 are successors (children) of node 12. Similarly, we define the ancestors and descendants of a node (see Section 2.2 for these definitions). For example, in Figure 8.9(a) node 6 has nodes 1, 2, 3, 4, 5, 6 as its descendants, and nodes 2, 5, and 6 are the ancestors of node 2. Notice that, according to our definitions, each node is its own ancestor and descendant.

This data structure supports the following six operations:

find-root(i). Find and return the root of the tree containing node i .

find-value(i). Find and return the value of the tree arc leaving node i . If i is a root node, return the value ∞ .

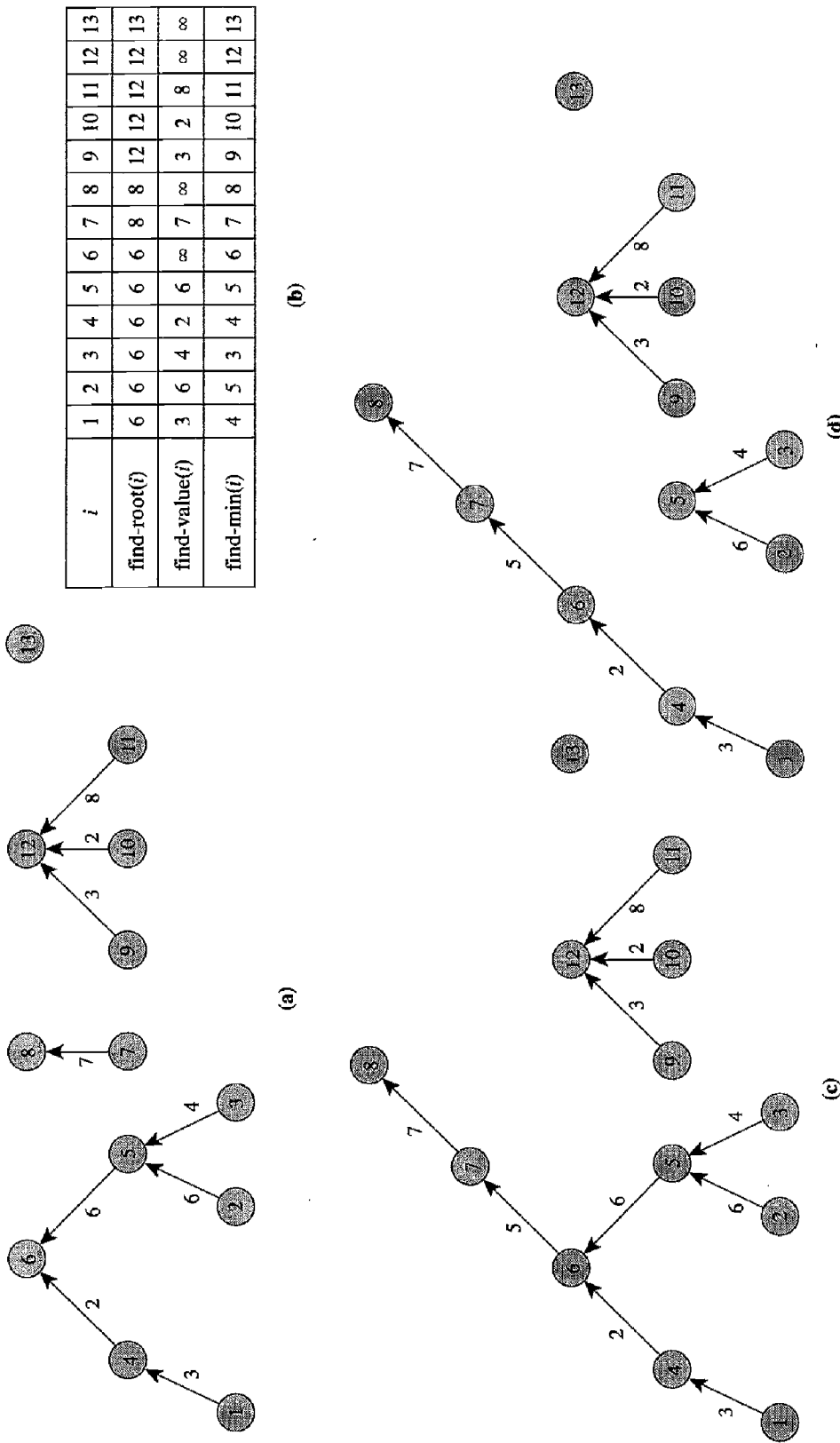


Figure 8.9 Illustrating various operations on dynamic trees: (a) collection of rooted trees; (b) results of the find-root, find-value, and find-min operations; (c) rooted trees after performing link (6, 7, 5); (d) rooted trees after performing cut (5).

find-min(i). Find and return the ancestor w of i with the minimum value of *find-value(w)*. In case of a tie, chose the node w closest to the tree root.

Figure 8.9(b) shows the results of the operations *find-root(i)*, *find-value(i)*, and *find-min(i)* performed for different nodes i .

change-value(i, val). Add a real number val to the value of every arc along the path from node i to *find-root(i)*. For example, if we execute *change-value(1, 3)* for the dynamic tree shown in Figure 8.9(a), we add 3 to the values of arcs (1, 4) and (4, 6) and these values become 6 and 5, respectively.

link(i, j, val). This operation assumes that i is a tree root and that i and j belong to different trees. The operation combines the trees containing nodes i and j by making node j the parent of node i and giving arc (i, j) the value val . As an illustration, if we perform the operation *link(6, 7, 5)* on our example, we obtain the trees shown in Figure 8.9(c).

cut(i). Break the tree containing node i into two trees by deleting the arc joining node i to its parent and returning the value of the deleted arc. We perform this operation when i is not a tree root. For example, if we execute *cut(5)* on trees in Figure 8.9(c), we delete arc (5, 6) and return its value 6. Figure 8.9(d) gives the new collection of trees.

The following important result, which we state without proof, lies at the heart of the efficiency of the dynamic tree data structure.

Lemma 8.9. *If z is the maximum tree size (i.e., maximum number of nodes in any tree), a sequence of l tree operations, starting with an initial collection of singleton trees, requires a total of $O(l \log(z + l))$ time.* ♦

The dynamic tree implementation stores the values of tree arcs only implicitly. If we were to store these values explicitly, the operation *change-value* on a tree of size z might require $O(z)$ time (if this tree happens to be a path), which is computationally excessive for most applications. Storing the values implicitly allows us to update the values in only $O(\log z)$ time. How the values are actually stored and manipulated is beyond the scope of this book.

How might we use the dynamic tree data structure to improve the computational performance of network flow algorithms. Let us use the shortest augmenting path algorithm as an illustration. The following basic idea underlies the algorithmic speed-up. In the dynamic tree implementation, each arc in the rooted tree is an admissible arc [recall that an arc (i, j) is admissible if $r_j > 0$ and $d(i) = d(j) + 1$]. The value of an arc is its residual capacity. For example, consider the residual network given in Figure 8.10(a), which shows the distance labels next to the nodes and residual capacities next to the arcs. Observe that in this network, every arc, except the arc (12, 13), is admissible; moreover, the residual capacity of every arc is 2, except for the arc (12, 14) whose residual capacity is 1. Figure 8.10(b) shows one collection of rooted trees for this example. Notice that although every tree arc is admissible, every admissible arc need not be in some tree. Consequently, for a given set of admissible arcs, many collections of rooted trees are possible.

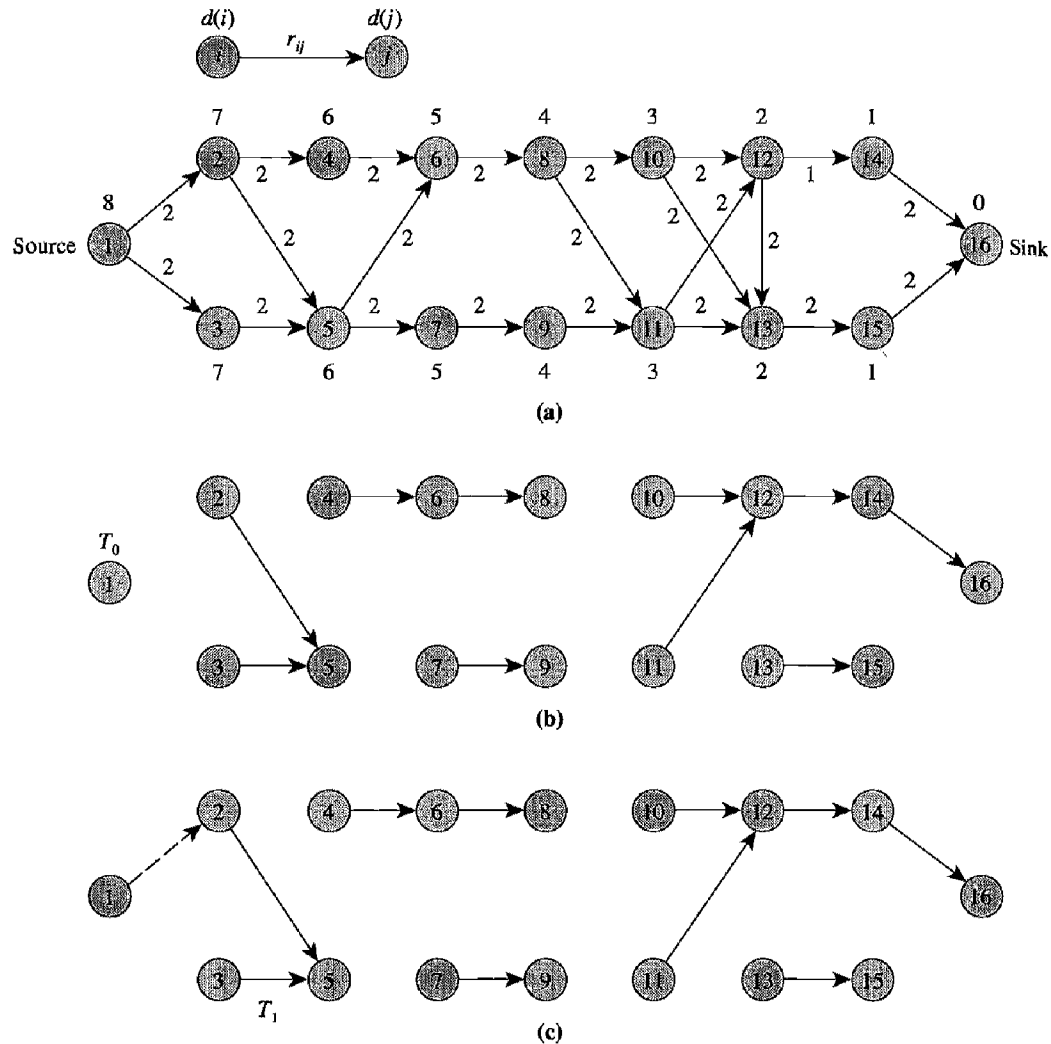


Figure 8.10 Finding an augmenting path in the dynamic tree implementations.

Before describing this implementation formally, we first show how the algorithm works on our example. It maintains a rooted tree containing the source node and progressively expands this tree until it contains the sink node, at which point the algorithm performs an augmentation. To grow the tree containing the source, the algorithm repeatedly performs link operations. In our example, the algorithm starts with the singleton tree T_0 containing only the source node 1 [see Figure 8.10(b)]. It identifies an admissible arc emanating from node 1. Suppose that we select arc (1, 2). The algorithm performs the operation $\text{link}(1, 2, 2)$, which joins two rooted trees, giving us a larger tree T_1 containing node 1 [see Figure 8.10(c)]. The algorithm then identifies the root of T_1 , by performing the operation $\text{find-root}(1)$, which identifies node 5. The algorithm tries to find an admissible arc emanating from node 5. Suppose that the algorithm selects the arc (5, 6). The algorithm performs the operation $\text{link}(5, 6, 2)$ and obtains a larger tree T_2 containing node 1 [see Figure 8.10(d)]. In the next iteration, the algorithm identifies node 8 as the root of T_2 . Suppose that

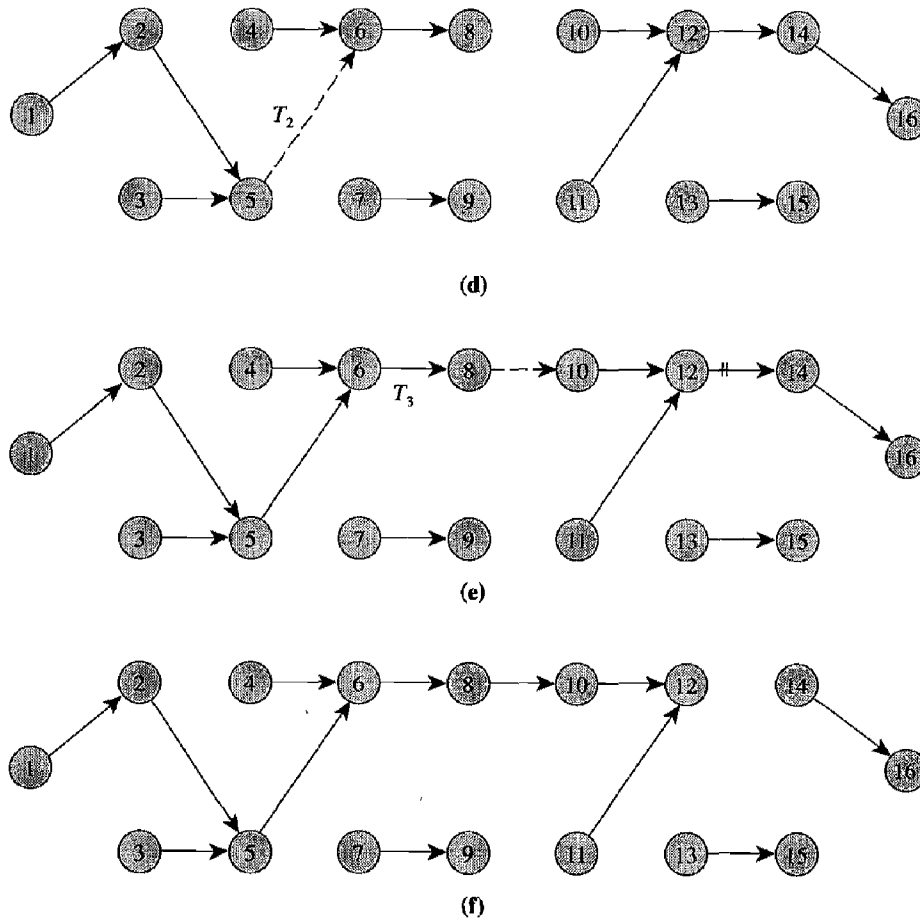


Figure 8.10 (Continued)

the algorithm selects arc $(8, 10)$ as the admissible arc emanating from node 8. The algorithm performs the operation $\text{link}(8, 10, 2)$ and obtains a rooted tree T_3 that contains both the source and sink nodes [see Figure 8.10(e)].

Observe that the unique path from the source to the sink in T_3 is an admissible path since by construction every arc in a rooted tree is admissible. The residual capacity of this path is the minimum value of the arcs in this path. How can we determine this value? Recall that the operation $\text{find-min}(1)$ would determine an ancestor of node 1 with the minimum value of find-value , which is node 12 in our example. Performing $\text{find-value}(12)$ will give us the residual capacity of this path, which is 1 in this case. We have thus discovered the possibility of augmenting 1 unit of flow along the admissible path and that arc $(12, 14)$ is the blocking arc. We perform the augmentation by executing $\text{change-value}(1, -1)$. This augmentation reduces the residual capacity of arc $(12, 14)$ to zero. The arc $(12, 14)$ now becomes inadmissible and we must drop it from the collection of rooted trees. We do so by performing $\text{cut}(12)$. This operation gives us the collection of rooted trees shown in Figure 8.10(f).

To better understand other situations that might occur, let us execute the algorithm for one more iteration. We apply the dynamic tree algorithm starting with the collection of rooted trees given in Figure 8.10(f). Node 12 is the root of the tree

containing node 1. But node 12 has no outgoing admissible arc; so we relabel node 12. This relabeling increases the distance label of node 12 to 3. Consequently, arcs (10, 12) and (11, 12) become inadmissible and we must drop them from the collection of rooted trees. We do so by performing $\text{cut}(10)$ and $\text{cut}(11)$, giving us the rooted trees shown in Figure 8.11(a). The algorithm again executes $\text{find-root}(1)$ and finds node 10 as the root of the tree containing node 1. In the next two operations, the algorithm adds arcs (10, 13) and (15, 16); Figure 8.11(b) and 8.11(c) shows the corresponding trees.

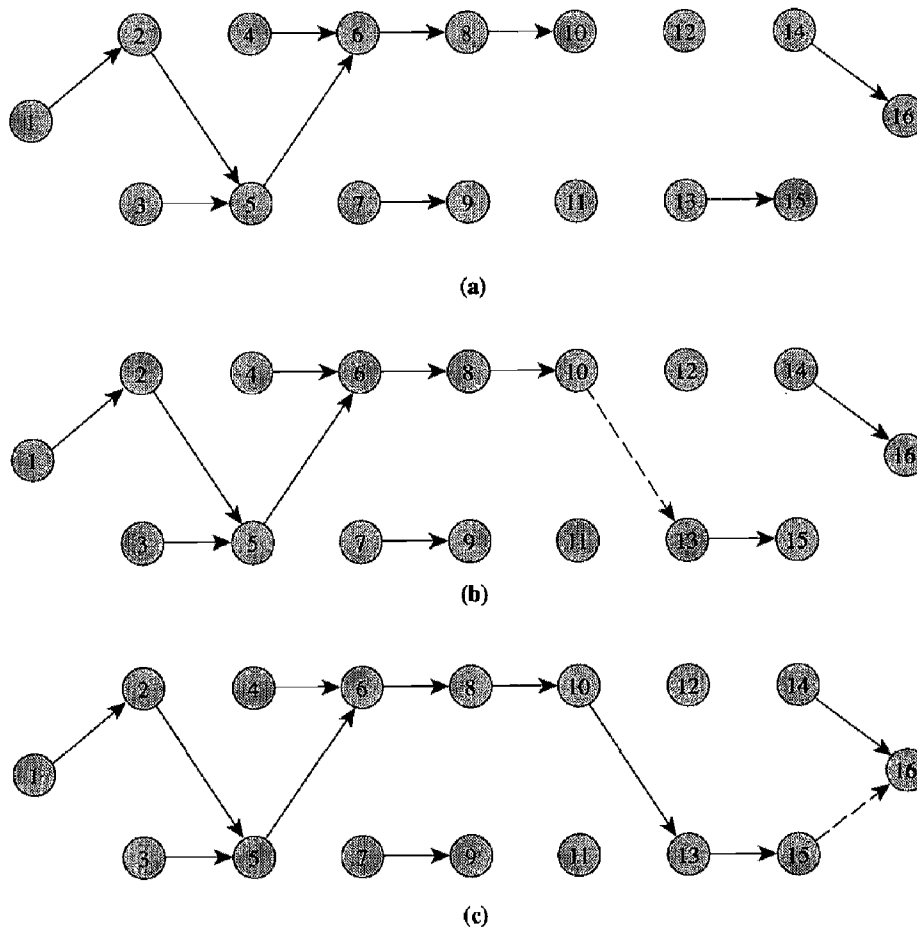


Figure 8.11 Another augmentation using dynamic trees.

Figures 8.12 and 8.13 give a formal statement of the algorithm. After offering explanatory comments, we consider a worst-case analysis of the algorithm. The algorithm is same as the one we presented in Section 7.4 except that it performs the procedures advance, retreat, and augment differently using trees. The first two procedures, tree-advance and tree-retreat, are straightforward, but the tree-augment procedure requires some explanation. If node p is an ancestor of node s with the minimum value of $\text{find-value}(p)$ and, among such nodes in the path, it is closest to the sink, then $\text{find-value}(p)$ gives the residual capacity of the augmenting path. The operation $\text{change}(s, -\delta)$ implicitly updates the residual capacities of all the arcs in

```

algorithm tree-augmenting-path;
begin
   $x := 0$ ;
  perform a reverse breadth-first search of the residual network
    from node  $t$  to obtain the distance labels  $d(l)$ ;
  let  $T$  be the collection of all singleton nodes;
   $i := s$ ;
  while  $d(s) < n$  do
    begin
      if  $i$  has an admissible arc then tree-advance( $i$ )
      else tree-retreat( $i$ );
      if  $i = t$  then tree-augment;
    end;
  end;

```

Figure 8.12 Dynamic tree implementation of the shortest augmenting path algorithm.

the augmenting path. This augmentation might cause the capacity of more than one arc in the path to become zero. The **while** loop identifies all such arcs, one by one, and deletes them from the collection of rooted trees.

We now consider the worst-case complexity of the algorithm. Why is the dynamic tree implementation more efficient than the original implementation of the shortest augmenting path algorithm? The bottleneck operations in the original shortest augmenting path algorithm are the advance and augment operations, which require $O(n^2m)$ time. Each advance operation in the original algorithm adds one arc;

```

procedure tree-advance( $i$ );
begin
  let  $(i, j)$  be an admissible arc in  $A(i)$ ;
  link( $i, j, r_{ij}$ );
   $i := \text{find-root}(j)$ ;
end;

```

(a)

```

procedure tree-retreat( $i$ );
begin
   $d(i) := \min\{d(j) + 1 : (i, j) \in A(i) \text{ and } r_{ij} > 0\}$ ;
  for each tree arc  $(k, i)$  do cut( $k$ );
   $i := \text{find-root}(s)$ ;
end;

```

(b)

```

procedure tree-augment;
begin
   $p := \text{find-min}(s)$ ;
   $\delta := \text{find-value}(p)$ ;
  change-value( $s, -\delta$ );
  while  $\text{find-value}(p) = 0$  do cut( $p$ ) and set  $p := \text{find-min}(s)$ ;
   $i := \text{find-root}(s)$ ;
end;

```

(c)

Figure 8.13 Procedures of the tree-augmenting-path algorithm.

in contrast, the tree implementation adds a collection of arcs using the link operation. Thus the dynamic tree implementation substantially reduces the number of executions of the link operation. Similarly, while augmenting flow, the tree implementation augments flow over a collection of arcs by performing the operation `change-value`, thus again substantially reducing the number of required updates.

We now obtain a bound on the number of times the algorithm performs various tree operations. We will show that the algorithm performs each of the tree operations $O(nm)$ times. In deriving these bounds, we make use of the results of Lemma 7.9, proved in Section 7.4.

cut(j). The algorithm performs this operation during the tree-retreat and tree-augment operations. During the `tree-retreat(i)` operation, the algorithm might perform this operation as many times as the number of incoming arcs at node i . Since this operation relabels node i , and we can relabel a node at most n times, these operations sum to $O(n^2)$ over all nodes. Furthermore, during the tree-augment operation, we perform the cut operation for each arc saturated during an augmentation. Since the total number of arc saturations is $O(nm)$, the number of these operations sums to $O(nm)$.

link(i, j, val). Each link operation adds an arc to the collection of rooted trees. Observe that if an arc enters a rooted tree, it remains there until a cut operation deletes it from the tree. Therefore, the number of link operations is at most $(n - 1)$ plus the number of cut operations. The term $(n - 1)$ arises because initially the collection might contain no arc, and finally, it might contain as many as $(n - 1)$ arcs. Consequently, the total number of link operations is also $O(nm)$. Since each tree-advance operation performs a link operation, the previous result also implies an $O(nm)$ bound on the total number of tree-advance operations.

change-value(i, val). The algorithm performs this operation once per augmentation. Since the number of augmentations is at most $nm/2$, we immediately obtain a bound of $O(nm)$ on the number of change-value operations.

find-value(i) and find-min(i). The algorithm performs each of these two operations once per augmentation and once for each arc saturated during the augmentation. These observations imply a bound of $O(nm)$ on the number of executions of these two operations.

find-root(i). The algorithm performs this operation once during each execution of the tree-advance, tree-augment, and tree-retreat operations. Since the algorithm executes the first two operations $O(nm)$ times and the third operation $O(n^2)$ times, it executes the find-root operation $O(nm)$ times.

Using simple arguments, we have now shown that the algorithm performs each of the six tree operations $O(nm)$ times. It performs each tree operation on a tree of maximum size n . The use of Lemma 8.9 establishes the following important result.

Theorem 8.10. *The dynamic tree implementation of the shortest augmenting path algorithm solves the maximum flow problem in $O(nm \log n)$ time.* ♦

Although this result establishes the theoretical utility of the dynamic tree data structure for improving the worst-case complexity of the shortest augmenting path algorithm, the practical value of this data structure is doubtful. The dynamic tree implementation reduces the time for performing advance and augment operations from $O(n^2m)$ to $O(nm \log n)$, but simultaneously increases the time of performing retreat operations from $O(nm)$ to $O(nm \log n)$. Empirical experience shows that the retreat operation is one of the bottleneck operations in practice. Since the dynamic tree data structure increases the running time of a bottleneck operation, the use of this data structure actually slows down the algorithm in practice. Furthermore, this data structure introduces substantial overhead (i.e., a large-constant factor of work is associated with each tree operation), thus making it of limited practical utility.

8.6 NETWORK CONNECTIVITY

The connectivity of a network is an important measure of the network's reliability or stability. The *arc connectivity* of a connected network is the minimum number of arcs whose removal from the network disconnects it into two or more components. In this section we suggest algorithms for solving the arc connectivity problem on undirected networks. The algorithms for solving connectivity problems on directed networks are different from those we will be discussing; we consider these algorithms in the exercises for this chapter. To conform with this choice of coverage, in this section by the term "network" we will invariably mean an undirected (connected) network.

The *node connectivity* of a network equals the minimum number of nodes whose deletion from the network disconnects it into two or more components. We discuss issues related to node connectivity in Exercise 8.35. We begin by defining several terms. A *disconnecting set* is a set of arcs whose deletion from the network disconnects it into two or more components. Therefore, the arc connectivity of a network equals the minimum cardinality of any disconnecting set; we refer to this set of arcs as a *minimum disconnecting set*.

The arc connectivity of a pair of nodes i and j is the minimum number of arcs whose removal from the network disconnects these two nodes. We represent the pair of nodes i and j by $[i, j]$ and the arc connectivity of this pair by $\alpha[i, j]$. We also let $\alpha(G)$ denote the arc connectivity of the network G . Consequently,

$$\alpha(G) = \min\{\alpha[i, j] : [i, j] \in N \times N\}. \quad (8.10)$$

We first bring together some elementary facts concerning the arc connectivity of a network; we ask the reader to prove these properties in Exercise 8.29.

Property 8.11.

- (a) $\alpha[i, j] = \alpha[j, i]$ for every pair of nodes $[i, j]$.
- (b) The arc connectivity of a network cannot exceed the minimum degree of nodes in the network. Therefore, $\alpha(G) \leq \lfloor m/n \rfloor$.

- (c) Any minimum disconnecting set partitions the network into exactly two components.
- (d) The arc connectivity of a spanning tree equals 1.
- (e) The arc connectivity of a cycle equals 2.

Let δ denote the minimum degree of a node in the network and let node p be a node with degree equal to δ . Property 8.11(b) implies that $\alpha(G) \leq \delta \leq \lfloor m/n \rfloor$. Since a minimum disconnecting set partitions the node set into exactly two components $S^* \subset N$ and $\bar{S}^* = N - S^*$, we can represent this cut by the notation $[S^*, \bar{S}^*]$. We assume, without any loss of generality, that node $p \in S^*$.

Our development in Chapter 6 provides us with a means for determining the arc connectivity of a network. Theorem 6.7 states that the minimum number of arcs in a network whose removal disconnects a specified pair of source and sink nodes equals the maximum number of arc-disjoint paths from the source to the sink. Furthermore, the proof of this theorem shows that we can obtain the maximum number of arc-disjoint paths from the source to the sink by solving a maximum flow problem in a network G whose arcs all have capacity equal to 1. Thus, to determine the arc connectivity of a network, we need to solve a unit capacity maximum flow problem between every pair of nodes (by varying the source and sink nodes); the minimum value among such flows is $\alpha(G)$. Since solving a unit capacity maximum flow problem requires $O(\min\{n^{2/3}m, m^{3/2}\})$ time (see Section 8.2), this approach produces an algorithm running in time $O(n^2 \cdot \min\{n^{2/3}m, m^{3/2}\})$.

We can easily improve on this approach by a factor of n using the following idea. Consider a node $k \in \bar{S}^*$ and recall that node $p \in S^*$. Since the cut $[S^*, \bar{S}^*]$ disconnects nodes p and k , the minimum cardinality of a set of arcs that will disconnect these two nodes is at most $|[S^*, \bar{S}^*]|$. That is,

$$\alpha[p, k] \leq |[S^*, \bar{S}^*]|. \quad (8.11)$$

Next observe that $[S^*, \bar{S}^*]$ is a minimum disconnecting set of the network. The definition (8.10) of $\alpha(G)$ implies that

$$\alpha[p, k] \geq |[S^*, \bar{S}^*]|. \quad (8.12)$$

Using (8.11) and (8.12), we see that

$$\alpha[p, k] = |[S^*, \bar{S}^*]|. \quad (8.13)$$

The preceding observations imply that if we compute $\alpha[p, j]$ for all j , the minimum among these numbers equals $\alpha(G)$. To summarize the discussion, we can write

$$\alpha(G) = \min\{\alpha[p, j] : j \in N - \{p\}\}.$$

The preceding approach permits us to determine the arc connectivity of a network by solving $(n - 1)$ unit capacity maximum flow problems, requiring $O(\min\{n^{5/3}m, nm^{3/2}\})$ time. We can improve this time bound for sparse networks by solving these maximum flow problems using the labeling algorithm described in Section 6.5 instead of the specialized unit capacity algorithms described in Section 8.2. The labeling algorithm will perform at most $\lfloor m/n \rfloor$ augmentations to solve each maximum flow problem (because the degree of node p is $\delta \leq \lfloor m/n \rfloor$) and would require $O(m^2/n)$ time. This approach requires $O(m^2)$ time to solve all the maximum

flow problems. Since $nm^{3/2} \geq m^2$, we can determine the arc connectivity of a network in $O(\min\{n^{5/3}m, m^2\})$ time. This algorithm is by no means the best algorithm for determining the arc connectivity of a network. We next describe an algorithm that computes arc connectivity in only $O(nm)$ time.

Just as the preceding algorithm determines the minimum cardinality of a set of arcs between pairs of nodes, the improved algorithm determines the minimum cardinality of a set of arcs that disconnects every node in a set S from some node $k \in \bar{S}$; we denote this number by $\alpha[S, k]$. We can compute $\alpha[S, k]$ using the labeling algorithm as follows: We allow the augmentation to start at any node in S but end only at node k . When the labeling algorithm terminates, the network contains no directed path from any node in S to node k . At this point the set of labeled nodes defines a cut in the network and the number of forward arcs in the cut is $\alpha[S, k]$.

Our preceding algorithm determines the arc connectivity of a network by computing $\alpha[p, j]$ for each node $j \in N - \{p\}$ and taking the minimum of these numbers. The correctness of this approach uses the fact that $\alpha(G)$ equals $\alpha[p, j]$ for some choice of the nodes p and j . Our improved algorithm determines the arc connectivity of a network by computing $\alpha[S, k]$ for at most $(n - 1)$ combinations of S and k and taking the minimum of these numbers. The algorithm selects the combinations S and k quite cleverly so that (1) for at least one combination of S and k , $\alpha[S, k] = \alpha(G)$; and (2) the labeling algorithm can compute $\alpha[S, k]$ for every combination in an average of $O(m)$ time because most augmentations involve only two arcs. Therefore this algorithm determines the arc connectivity of a network in $O(nm)$ time.

Before describing the algorithm, we first introduce some notation. For any set S of nodes, we let $\text{neighbor}(S)$ denote the set of nodes in \bar{S} that are adjacent to some node in S , and $\text{nonneighbor}(S)$ as the set of nodes in \bar{S} that are not adjacent to any node in S . Consequently, $N = S \cup \text{neighbor}(S) \cup \text{nonneighbor}(S)$. Our improved arc connectivity algorithm depends on the following crucial result.

Lemma 8.12. *Let δ be the minimum node degree of a network G and let $[S^*, \bar{S}^*]$ denote a minimum disconnecting set of the network. Suppose that $\alpha(G) \leq \delta - 1$. Then for any set $S \subseteq S^*$, $\text{nonneighbor}(S)$ is nonempty.*

Proof. We first notice that the maximum number of arcs emanating from nodes in \bar{S}^* is $|\bar{S}^*|(|\bar{S}^*| - 1) + \alpha(G)$ because any such arc either has both its endpoints in \bar{S}^* or belongs to the minimum disconnecting set. Next notice that the minimum number of arcs emanating from the nodes in \bar{S}^* is $\delta|\bar{S}^*|$ because δ is the minimum node degree. Therefore,

$$|\bar{S}^*|(|\bar{S}^*| - 1) + \alpha(G) \geq |\bar{S}^*|\delta.$$

Adding δ to both the sides of this inequality and simplifying the expression gives

$$(|\bar{S}^*| - 1)(|\bar{S}^*| - \delta) \geq \delta - \alpha(G) \geq 1.$$

The last inequality in this expression follows from the fact $\alpha(G) \leq \delta - 1$. Notice that the inequality $(|\bar{S}^*| - 1)(|\bar{S}^*| - \delta) \geq 1$ implies that both the terms to the left are at least one. Thus $|\bar{S}^*| \geq \delta + 1$; that is, the set \bar{S}^* contains at least $\delta + 1$ nodes. Since the cut $[S^*, \bar{S}^*]$ contains fewer than δ arcs, at least one of the nodes in \bar{S}^* is not adjacent to any node in S^* . Consequently, the set $\text{nonneighbor}(S)$ must be nonempty, which establishes the lemma. ♦

The improved arc connectivity algorithm works as follows. It starts with $S = \{p\}$, selects a node $k \in \text{nonneighbor}(S)$, and computes $\alpha[S, k]$. It then adds node k to S , updates the sets $\text{neighbor}(S)$ and $\text{nonneighbor}(S)$, selects another node $k \in \text{nonneighbor}(S)$ and computes $\alpha[S, k]$. It repeats this operation until the set $\text{nonneighbor}(S)$ is empty. The minimum value of $\alpha[S, k]$, obtained over all the iterations, is $\alpha(G)$. Figure 8.14 gives a formal description of this algorithm.

```

algorithm arc connectivity;
begin
  let  $p$  be a minimum degree node in the network and  $\delta$  be its degree;
  set  $S^* := \{p\}$  and  $\alpha^* := \delta$ ;
  set  $S := \{p\}$ ;
  initialize  $\text{neighbor}(S)$  and  $\text{nonneighbor}(S)$ ;
  while  $\text{nonneighbor}(S)$  is nonempty do
    begin
      select a node  $k \in \text{nonneighbor}(S)$ ;
      compute  $\alpha[S, k]$  using the labeling algorithm for the maximum flow
        problem and let  $[R, \bar{R}]$  be the corresponding disconnecting cut;
      if  $\alpha^* > \alpha[S, k]$  then set  $\alpha^* := \alpha[S, k]$  and  $[S^*, \bar{S}^*] := [R, \bar{R}]$ ;
      add node  $k$  to  $S$  and update  $\text{neighbor}(S)$ ,  $\text{nonneighbor}(S)$ ;
    end;
  end;

```

Figure 8.14 Arc connectivity algorithm.

To establish the correctness of the arc connectivity algorithm, let $[S^*, \bar{S}^*]$ denote the minimum disconnecting set. We consider two cases: when $\alpha(G) = \delta$ and when $\alpha(G) \leq \delta - 1$. If $\alpha(G) = \delta$, the algorithmic description in Figure 8.14 implies that the algorithm would terminate with $[p, N - \{p\}]$ as the minimum disconnecting set. Now suppose $\alpha(G) \leq \delta - 1$. During its execution, the arc connectivity algorithm determines $\alpha[S, k]$ for different combinations of S and k ; we need to show that at some iteration, $\alpha[S, k]$ would equal $\alpha(G)$. We establish this result by proving that at some iteration, $S \subseteq S^*$ and $k \in \bar{S}^*$, in which case $\alpha[S, k] = \alpha(G)$ because the cut $[S^*, \bar{S}^*]$ disconnects every node in S from node k . Notice that initially S^* contains S (because both start with p as their only element), and finally it does not because, from Lemma 8.9, as long as S^* contains S , $\text{nonneighbor}(S)$ is nonempty and the algorithm can add nodes to S . Now consider the last iteration for which $S \subseteq S^*$. At this iteration, the algorithm selects a node k that must be in \bar{S}^* because $S \cup \{k\} \not\subseteq S^*$. But then $\alpha[S, k] = \alpha(G)$ because the cut $[S^*, \bar{S}^*]$ disconnects S from node k . This conclusion shows that the arc connectivity algorithm correctly solves the connectivity problem.

We next analyze the complexity of the arc connectivity algorithm. The algorithm uses the labeling algorithm described in Section 6.5 to compute $\alpha[S, k]$; suppose that the labeling algorithm examines labeled nodes in the first-in, first-out order so that it augments flow along shortest paths in the residual network. Each augmenting path starts at a node in S , terminates at the node k , and is one of two types: Its last internal node is in $\text{neighbor}(S)$ or it is in $\text{nonneighbor}(S)$. All augmenting paths of the first type are of length 2; within an iteration, we can find any such path in a total of $O(n)$ time (why?), so augmentations of the first type require a total of $O(n^2)$ time in all iterations. Detecting an augmenting path of the second type requires

$O(m)$ time; it is possible to show, however, that in all the applications of the labeling algorithm in various iterations, we never encounter more than n such augmenting paths. To see this, consider an augmenting path of the second type which contains node $l \in \text{nonneighbor}(S)$ as the last internal node in the path. At the end of this iteration, the algorithm will add node k to S ; as a result, it adds node l to $\text{neighbor}(S)$; the node will stay there until the algorithm terminates. So each time the algorithm performs an augmentation of the second type, it moves a node from the set $\text{nonneighbor}(S)$ to $\text{neighbor}(S)$. Consequently, the algorithm performs at most n augmentations of the second type and the total time for these augmentations will be $O(nm)$. The following theorem summarizes this discussion.

Theorem 8.13. *In $O(nm)$ time the arc connectivity algorithm correctly determines the arc connectivity of a network.* ♦

8.7 ALL-PAIRS MINIMUM VALUE CUT PROBLEM

In this section we study the all-pairs minimum value cut problem in undirected network, which is defined in the following manner. For a specific pair of nodes i and j , we define an $[i, j]$ cut as a set of arcs whose deletion from the network disconnects the network into two components S_{ij} and \bar{S}_{ij} so that nodes i and j belong to different components (i.e., if $i \in S_{ij}$, then $j \in \bar{S}_{ij}$; and if $i \in \bar{S}_{ij}$, then $j \in S_{ij}$). We refer to this $[i, j]$ cut as $[S_{ij}, \bar{S}_{ij}]$ and say that this cut *separates* nodes i and j . We associate with a cut $[S_{ij}, \bar{S}_{ij}]$, a *value* that is a function of arcs in the cut. A *minimum $[i, j]$ cut* is a cut whose value is minimum among all $[i, j]$ cuts. We let $[S_{ij}^*, \bar{S}_{ij}^*]$ denote a minimum value $[i, j]$ cut and let $v[i, j]$ denote its value. The all-pairs minimum value cut problem requires us to determine for all pairs of nodes i and j , a minimum value $[i, j]$ cut $[S_{ij}^*, \bar{S}_{ij}^*]$ and its value $v[i, j]$.

The definition of a cut implies that if $[S_{ij}, \bar{S}_{ij}]$ is an $[i, j]$ cut, it is also a $[j, i]$ cut. Therefore, $v[i, j] = v[j, i]$ for all pairs i and j of nodes. This observation implies that we can solve the all-pairs minimum value cut problem by invoking $n(n - 1)/2$ applications of any algorithm for the single pair minimum value cut problem. We can, however, do better. In this section we show that we can solve the all-pairs minimum value cut problem by invoking only $(n - 1)$ applications of the single-pair minimum value cut problem.

We first mention some specializations of the all-pairs minimum value cut problem on undirected networks. If we define the value of a cut as its capacity (i.e., the sum of capacities of arcs in the cut), the all-pairs minimum value cut problem would identify minimum cuts (as defined in Chapter 6) between all pairs of nodes. Since the minimum cut capacity equals the maximum flow value, we also obtain the maximum flow values between all pairs of nodes. Several other functions defined on a cut $[S_{ij}, \bar{S}_{ij}]$ are plausible, including (1) the number of arcs in the cut, (2) the capacity of the cut divided by the number of arcs in the cut, and (3) the capacity of the cut divided by $|S_{ij}| + |\bar{S}_{ij}|$.

We first state and prove an elementary lemma concerning minimum value cuts.

Lemma 8.14. *Let i_1, i_2, \dots, i_k be an (arbitrary) sequence of nodes. Then $v[i_1, i_k] \geq \min\{v[i_1, i_2], v[i_2, i_3], \dots, v[i_{k-1}, i_k]\}$.*

Proof. Let $i = i_1, j = i_k$, and $[S_{ij}^*, \bar{S}_{ij}^*]$ be the minimum value $[i, j]$ cut. Consider the sequence of nodes i_1, i_2, \dots, i_k in order and identify the smallest index r satisfying the property that i_r and i_{r+1} are in different components of the cut $[S_{ij}^*, \bar{S}_{ij}^*]$. Such an index must exist because, by definition, $i_1 = i \in S_{ij}^*$ and $i_k = j \notin S_{ij}^*$. Therefore, $[S_{ij}^*, \bar{S}_{ij}^*]$ is also an $[i_r, i_{r+1}]$ cut, which implies that the value of the minimum value $[i_r, i_{r+1}]$ cut will be no more than the value of the cut $[S_{ij}^*, \bar{S}_{ij}^*]$. In other words,

$$v[i_1, i_k] \geq v[i_r, i_{r+1}] \geq \min\{v[i_1, i_2], v[i_2, i_3], \dots, v[i_{k-1}, i_k]\}, \quad (8.14)$$

which is the desired conclusion of the lemma. ♦

Lemma 8.14 has several interesting implications. Select any three nodes i, j , and k of the network and consider the minimum cut values $v[i, j]$, $v[j, k]$, and $v[k, i]$ between them. The inequality (8.14) implies that at least two of the values must be equal. For if these three values are distinct, then placing the smallest value on the left-hand side of (8.14) would contradict this inequality. Furthermore, it is possible to show that one of these values that is not equal to the other two must be the largest. Since for every three nodes, two of the three minimum cut values must be equal, it is conceivable that many of the $n(n - 1)/2$ cut values will be equal. Indeed, it is possible to show that the number of distinct minimum cut values is at most $(n - 1)$. This result is the subject of our next lemma. This lemma requires some background concerning the *maximum spanning tree problem* that we discuss in Chapter 13. In an undirected network G , with an associated value (or, profit) for each arc, the maximum spanning tree problem seeks a spanning tree T^* , from among all spanning trees, with the largest sum of the values of its arcs. In Theorem 13.4 we state the following optimality condition for the maximum spanning tree problem: A spanning tree T^* is a maximum spanning tree if and only if for every nontree arc (k, l) , the value of the arc (k, l) is less than or equal to the value of every arc in the unique tree path from node k to node l .

Lemma 8.15. *In the $n(n - 1)/2$ minimum cut values between all pairs of nodes, at most $(n - 1)$ values are distinct.*

Proof. We construct a *complete* undirected graph $G' = (N, A')$ with n nodes. We set the value of each arc $(i, j) \in A'$ equal to $v[i, j]$ and associate the cut $[S_{ij}^*, \bar{S}_{ij}^*]$ with this arc. Let T^* be a maximum spanning tree of G' . Clearly, $|T^*| = n - 1$. We shall prove that the value of every nontree arc is equal to the value of some tree arc in T^* and this result would imply the conclusion of the lemma.

Consider a nontree arc (k, l) of value $v[k, l]$. Let P denote the unique path in T^* between nodes k and l . The fact that T^* is a maximum spanning tree implies that the value of arc (k, l) is less than or equal to the value of every arc $(i, j) \in P$. Therefore,

$$v[k, l] \leq \min\{v[i, j] : (i, j) \in P\}. \quad (8.15)$$

Now consider the sequence of nodes in P that starts at node k and ends at node l . Lemma 8.14 implies that

$$v[k, l] \geq \min\{v[i, j] : (i, j) \in P\}. \quad (8.16)$$

The inequalities (8.15) and (8.16) together imply that

$$v[k, l] = \min[v[i, j] : (i, j) \in P].$$

Consequently, the value of arc (k, l) equals the minimum value of an arc in P , which completes the proof of the lemma. ♦

The preceding lemma implies that we can store the $n(n - 1)/2$ minimum cut values in the form of a spanning tree T^* with a cut value associated with every arc in the tree. To determine the minimum cut values $v[k, l]$ between a pair k and l of nodes, we simply traverse the tree path from node k to node l ; the cut value $v[k, l]$ equals the minimum value of any arc encountered in this path. Note that the preceding lemma only establishes the fact that there are at most $(n - 1)$ distinct minimum cut values and shows how to store them compactly in the form of a spanning tree. It does not, however, tell us whether we can determine these distinct cut values by solving $(n - 1)$ minimum cut problems, because the proof of the lemma requires the availability of minimum cut values between all node pairs which we do not have.

Now, we ask a related question. Just as we can concisely store the minimum cut values between all pairs of nodes by storing only $(n - 1)$ values, can we also store the minimum value cuts between all pairs of nodes concisely by storing only $(n - 1)$ cuts? Because the network has at most $(n - 1)$ distinct minimum cut values between $n(n - 1)/2$ pairs of nodes, does it have at most $(n - 1)$ distinct cuts that define the minimum cuts between all node pairs? In the following discussion we provide an affirmative answer to this question. Consider a pair k and l of nodes. Suppose that arc (i, j) is a minimum value arc in the path from node k to node l in T^* . Our preceding observations imply that $v[k, l] = v[i, j]$. Also notice that we have associated a cut $[S_{ij}^*, \bar{S}_{ij}^*]$ of value $v[i, j] = v[k, l]$ with the arc (i, j) ; this cut separates nodes i and j . Is $[S_{ij}^*, \bar{S}_{ij}^*]$ a minimum $[k, l]$ cut? It is if $[S_{ij}^*, \bar{S}_{ij}^*]$ separates nodes k and l , and it is not otherwise. If, indeed, $[S_{ij}^*, \bar{S}_{ij}^*]$ separates nodes k and l , and if the same result is true for every pair of nodes in the network, the cuts associated with arcs in T^* concisely store minimum value cuts between all pairs of nodes. We refer to such a tree T^* as a *separator tree*. In this section we show that every network G has a separator tree and that we can construct the separator tree by evaluating $(n - 1)$ single-pair minimum cut values. Before we describe this method, we restate the definition of the separator tree for easy future reference.

Separator tree. An undirected spanning tree T^* , with a minimum $[i, j]$ cut $[S_{ij}^*, \bar{S}_{ij}^*]$ of value $v[i, j]$ associated with each arc (i, j) , is a separator tree if it satisfies the following property for every nontree arc (k, l) : If arc (i, j) is the minimum value arc from node k to node l in T^* (breaking ties in a manner to be described later), $[S_{ij}^*, \bar{S}_{ij}^*]$ separates nodes k and l .

Our preceding discussion shows that we have reduced the all-pairs minimum value cut problem to a problem of obtaining a separator tree. Given the separator tree T^* , we determine the minimum $[k, l]$ cut as follows: We traverse the tree path from node k to node l ; the cut corresponding to the minimum value in the path (breaking ties appropriately) is a minimum $[k, l]$ cut.

We call a subtree of a separator tree a *separator subtree*. Our algorithm for

constructing a separator tree proceeds by constructing a separator subtree that spans an expanding set for nodes. It starts with a singleton node, adds one additional node to the separator subtree at every iteration, and terminates when the separator tree spans all the nodes. We add nodes to the separator subtree in the order $1, 2, 3, \dots, n$. Let T^{p-1} denote the separator subtree for the node set $\{1, 2, \dots, p-1\}$ and T^p denote the separator subtree for the node set $\{1, 2, \dots, p\}$. We obtain T^p from T^{p-1} by adding an arc, say (p, k) . The essential problem is to locate the node k incident to node p in T^p . Once we have located the node k , we identify a minimum value cut $[S_{pk}^*, \bar{S}_{pk}^*]$ between nodes p and k , and associate it with the arc (p, k) . We set the value of the arc (p, k) equal to $v[p, k]$.

As already mentioned, our algorithm for constructing the separator tree adds arcs to the separator subtree one by one. We associate an index, called an *order* index, with every arc in the separator subtree. The first arc added to the separator subtree has order index 1, the second arc added has order index 2, and so on. We use the order index to resolve ties while finding a minimum value arc between a pair of nodes. As a rule, whether we specify so or not in the subsequent discussion, we always resolve any tie in favor of the arc with the least order index (i.e., the arc that we added first to the separator subtree).

Figure 8.15 describes the procedure we use to locate the node $k \in T^{p-1}$ on which the arc (p, k) will be incident in T^p .

Note that in every iteration of the locate procedure, $T_\alpha \subseteq N_\alpha$ and $T_\beta \subseteq N_\beta$. This fact follows from the observation that for every $k \in T_\alpha$ and $l \in T_\beta$, the arc (α, β) is the minimum value arc in the path from node k to node l in the separator subtree T and, by its definition, the cut must separate node k and node l .

We illustrate the procedure locate using a numerical example. Consider a nine-node network with nodes numbered $1, 2, 3, \dots, 9$. Suppose that after five iterations, the separator subtree $T^{p-1} = T^6$ is as shown in Figure 8.16(a). The figure also shows the cut associated with each arc in the separator subtree (here we specify only $S_{\alpha\beta}^*$, because we can compute $\bar{S}_{\alpha\beta}^*$ by using $\bar{S}_{\alpha\beta}^* = N - S_{\alpha\beta}^*$). We next consider adding node 7 to the subtree. At this point, $T = T^6$ and the minimum value arc in T is $(4, 5)$. Examining S_{45}^* reveals that node 7 is on the same side of the cut as node

```

procedure locate( $T^{p-1}, p, k$ );
begin
     $T := T^{p-1}$ ;
    while  $T$  is not a singleton node do
        begin
            let  $(\alpha, \beta)$  be the minimum value arc in  $T$  (we break ties in favor of the arc with the
                smallest order index);
            let  $[S_{\alpha\beta}^*, \bar{S}_{\alpha\beta}^*]$  be the cut associated with the arc  $(\alpha, \beta)$ ;
            let the arc  $(\alpha, \beta)$  partition the tree  $T$  into the subtrees  $T_\alpha$ 
                and  $T_\beta$  so that  $\alpha \in T_\alpha$  and  $\beta \in T_\beta$ ;
            let the cut  $[S_{\alpha\beta}^*, \bar{S}_{\alpha\beta}^*]$  partition the node set  $N$  into the
                subsets  $N_\alpha$  and  $N_\beta$  so that  $\alpha \in N_\alpha$  and  $\beta \in N_\beta$ ;
            if  $p \in N_\alpha$  then set  $T := T_\alpha$  else set  $T := T_\beta$ ;
        end;
    set  $k$  equal to the singleton node in  $T$ ;
end;

```

Figure 8.15 Locate procedure.

4; therefore, we update T to be the subtree containing node 4. Figure 8.16(b) shows the tree T . Now arc $(2, 3)$ is the minimum value arc in T . Examining S_{23}^* reveals that node 7 is on the same side of the cut as node 2; so we update T so that it is the subtree containing node 2. At this point, the tree T is a singleton, node 2. We set $k = 2$ and terminate the procedure. We next add arc $(2, 7)$ to the separator subtree, obtain a minimum value cut between the nodes 7 and 2, and associate this cut with the arc $(2, 7)$. Let $v[7, 2] = 5$ and $S_{72}^* = \{7\}$. Figure 8.16(c) shows the separator subtree spanning the nodes 1 through 7.

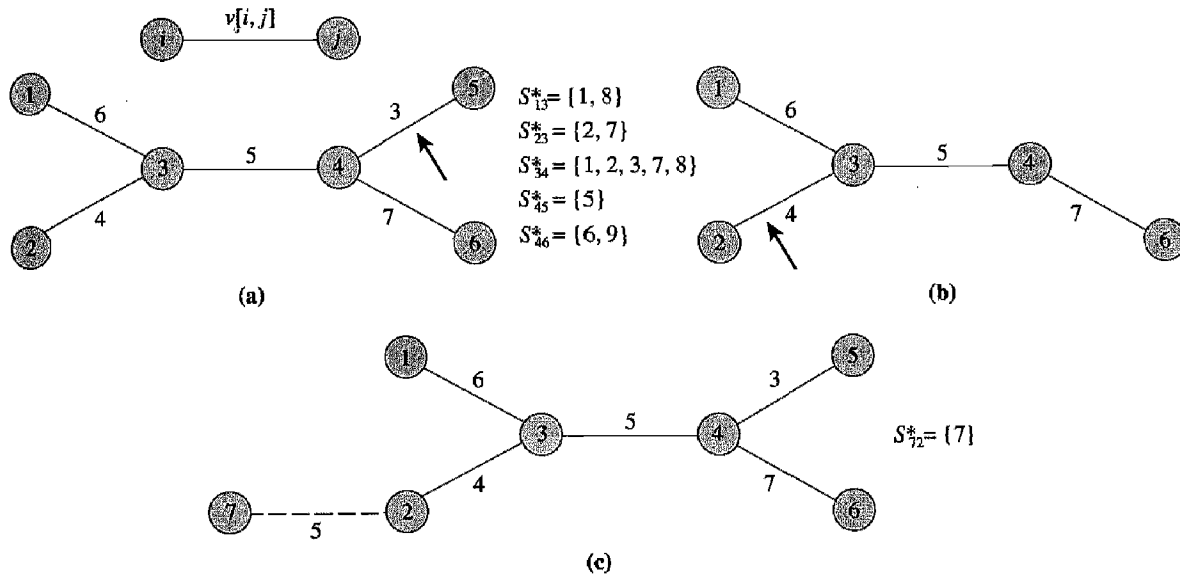


Figure 8.16 Illustrating the procedure locate.

We are now in a position to prove that the subtree T^p is a separator subtree spanning the nodes $\{1, 2, \dots, p\}$. Since, by our inductive hypothesis, T^{p-1} is a separator subtree on the nodes $\{1, 2, \dots, p-1\}$, our proof amounts to establishing the following result for every node $l \in \{1, 2, \dots, p-1\}$: If (i, j) is the minimum value arc in T^p in the path from node p to node l (with ties broken appropriately), the cut $[S_{ij}^*, \bar{S}_{ij}^*]$ separates the nodes p and l . We prove this result in the following lemma.

Lemma 8.16. *For any node $l \in T^{p-1}$, if (i, j) is the minimum value arc in T^p in the path from node p to node l (when we break ties in favor of the arc with the least order index), $[S_{ij}^*, \bar{S}_{ij}^*]$ separates nodes p and l .*

Proof. We consider two possibilities for the arc (i, j) .

Case 1: $(i, j) = (p, k)$. Let P denote the tree path in T^p from node k to node l . The situation $(i, j) = (p, k)$ can occur only when arc (p, k) is the unique minimum value arc in T^p in P , for otherwise, the tie will be broken in favor of an arc other than the arc (p, k) (why?). Thus

$$v[p, k] < v[g, h] \quad \text{for every arc } (g, h) \in P. \quad (8.17)$$

Next consider any arc $(g, h) \in P$. We claim that both the nodes g and h must belong to the same component of the cut $[S_{pk}^*, \bar{S}_{pk}^*]$; for otherwise, $[S_{pk}^*, \bar{S}_{pk}^*]$ will also separate nodes g and h , so $v[p, k] \geq v[g, h]$, contradicting (8.17). Using this argument inductively for all arcs in P , we see that all the nodes in the path P (that starts at node k and ends at node l) must belong to the same component of the cut $[S_{pk}^*, \bar{S}_{pk}^*]$. Since the cut $[S_{pk}^*, \bar{S}_{pk}^*]$ separates nodes p and k , it also separates nodes p and l .

Case 2: $(i, j) \neq (p, k)$. We examine this case using the locate procedure. At the beginning of the locate procedure, $T = T^p$, and at every iteration the size of the tree becomes smaller, until finally, $T = \{k\}$. Consider the iteration when T contains both the nodes k and l , but in the next iteration the tree does not contain node l . Let P denote the path in T from node k to node l in this iteration. It is easy to see that the arc (α, β) selected by the locate procedure in this iteration must belong to the path P . By definition, (α, β) is the minimum value arc in T , with ties broken appropriately. Since T contains the path P , (α, β) is also a minimum value arc in P . Now notice from the statement of the lemma that arc (i, j) is defined as the minimum value arc in P , and since we break the tie in precisely the same manner, $(i, j) = (\alpha, \beta)$.

We next show that the cut $[S_{\alpha\beta}^*, \bar{S}_{\alpha\beta}^*]$ separates node p and node l . We recommend that the reader refers to Figure 8.17 while reading the remainder of the proof. Consider the same iteration of the locate procedure considered in the preceding paragraph, and let T_α , N_α , T_β , N_β be defined as in Figure 8.15. We have observed previously that $T_\alpha \subseteq N_\alpha$ and $T_\beta \subseteq N_\beta$. We assume that $p \in N_\alpha$; a similar argument applies when $p \in N_\beta$. The procedure implies that when $p \in N_\alpha$, we set $T = T_\alpha$, implying that $k \in T_\alpha \subseteq N_\alpha$. Since the cut $[S_{\alpha\beta}^*, \bar{S}_{\alpha\beta}^*]$ separates node k from node l it also separates node p from node l . The proof of the lemma is complete. \blacklozenge

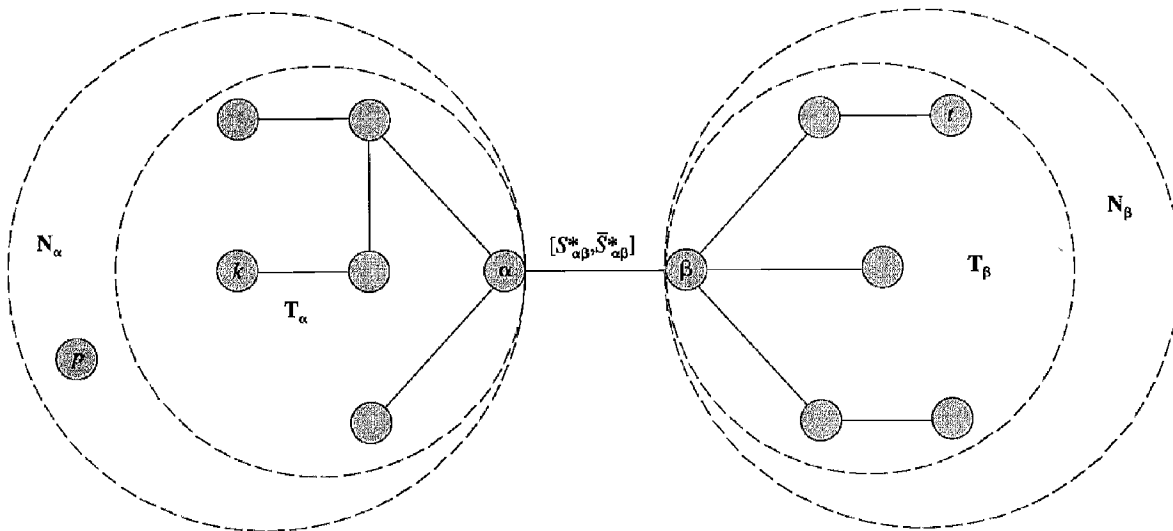


Figure 8.17 Proving Lemma 8.16.

Having proved the correctness of the all-pairs minimum cut algorithm, we next analyze its running time. The algorithm performs $(n - 1)$ iterations. In each iteration it executes the locate procedure and identifies the arc (p, k) to be added to the separator subtree. The reader can easily verify that an execution of the locate procedure requires $O(n^2)$ time. The algorithm then solves a minimum value cut problem and associates this cut and its value with the arc (p, k) . The following theorem is immediate.

Theorem 8.17. *Solving the all-pairs minimum value cut problem requires $O(n^3)$ time plus the time required to solve $(n - 1)$ instances of the single-pair minimum value cut problem.* ♦

To summarize, we have shown that the all-pairs minimum cut algorithm finds the minimum capacity cut separating node i and node j (i.e., with node i on one side of the cut and node j on the other side) for every node pair $[i, j]$. The max-flow min-cut theorem shows that this algorithm also determines the maximum flow values between every pair of nodes.

Suppose, instead, that we are given a directed graph. Let $f[i, j]$ denote the value of the minimum cut from node i to node j . We cannot use the all-pairs minimum cut algorithm to determine $f[i, j]$ for all node pairs $[i, j]$ for the following simple reason: The algorithm would determine the minimum value of a cut separating node i from node j , and this value is $\min\{f[i, j], f[j, i]\}$ because the minimum cut from node i to node j separates nodes i and j and so does the minimum cut from node j to node i . (We did not face this problem for undirected networks because the minimum cut from node i to node j is also a minimum cut from node j to node i .) If we let $v[i, j] = \min\{f[i, j], f[j, i]\}$, we can use the all-pairs minimum cut algorithm to determine $v[i, j]$ for each node pair $[i, j]$. Moreover, this algorithm relies on evaluating $v[i, j]$ for only $(n - 1)$ pairs of nodes. Since we can determine $v[i, j]$ by finding a maximum flow from node i to node j and from node j to node i , we can compute $v[i, j]$ for all node pairs by solving $(2n - 2)$ maximum flow problems.

We complete this section by describing an application of the all-pairs minimum value cut problem.

Application 8.3 Maximum and Minimum Arc Flows in a Feasible Flow

Consider the feasible flow problem that we discussed in Application 6.1. Assume that the network is uncapacitated and that it admits a feasible flow. For each arc $(i, j) \in A$, let α_{ij} denote the minimum arc flow that (i, j) can have in some feasible flow, and let β_{ij} denote the maximum arc flow that (i, j) can have in some feasible flow. We will show that we can determine α_{ij} for all node pairs $[i, j]$ by solving at most n maximum flow problems, and we can determine β_{ij} for all node pairs $[i, j]$ by an application of the all-pairs minimum value cut problem.

The problem of determining the maximum and minimum values of the arc flows in a feasible flow arises in the context of determining the statistical security of data (e.g., census data). Given a two-dimensional table A of size $p \times q$, suppose that we want to disclose the row sums r_i , the column sums c_j , and a subset of the matrix

elements. For security reasons (or to ensure confidentiality of the data), we would like to ensure that we have “disguised” the remaining matrix elements (or “hidden” entries). We wish to address the following question: Once we have disclosed the row and column sums and some matrix elements, how secure are the hidden entries? We address a related question: For each hidden element a_{ij} , what are the minimum and maximum values that a_{ij} can assume consistent with the data we have disclosed? If these two bounds are quite close, the element a_{ij} is not secure.

We assume that each revealed matrix element has value 0. We incur no loss of generality in making this assumption since we can replace a nonzero element a_{ij} by 0, replace r_i by $r_i - a_{ij}$, and replace c_j by $c_j - a_{ij}$. To conduct our analysis, we begin by constructing the bipartite network shown in Figure 8.18; in this network, each unrevealed matrix element a_{ij} corresponds to an arc from node i to node j . It is easy to see that every feasible flow in this network gives values of the matrix elements that are consistent with the row and column sums.

How might we compute the α_{ij} values? Let x^* be any feasible flow in the network (which we can determine by solving a maximum flow problem). In Section 11.2 we show how to convert each feasible flow into a feasible spanning tree solution; in this solution at most $(n - 1)$ arcs have positive flow. As a consequence, if x^* is a spanning tree solution, at least $(m - n + 1)$ arcs have zero flow; and therefore, $\alpha_{ij} = 0$ for each of these arcs. So we need to find the α_{ij} values for only the remaining $(n - 1)$ arcs. We claim that we can accomplish this task by solving at most $(n - 1)$ maximum flow problems. Let us consider a specific arc (i, j) . To determine the minimum flow on arc (i, j) , we find the maximum flow from node i to node j in $G(x^*)$ when we have deleted arc (i, j) . If the maximum flow from node i to node j has value k , we can reroute up to k units of flow from node i to node j and reduce the flow on arc (i, j) by the same amount. As a result, $\alpha_{ij} = \max\{0, x_{ij}^* - k\}$.

To determine the maximum possible flow on arc (i, j) , we determine the maximum flow from node j to node i in $G(x^*)$. If we can send k units from node j to node i in $G(x^*)$, then $\beta_{ij} = k$. To establish this result, suppose that the k units consist of x_{ij}^* units on arc (j, i) [which is the reversal of arc (i, j)], and $k - x_{ij}^*$ units that do not use arc (i, j) . Then, to determine the maximum flow on the arc (i, j) , we can

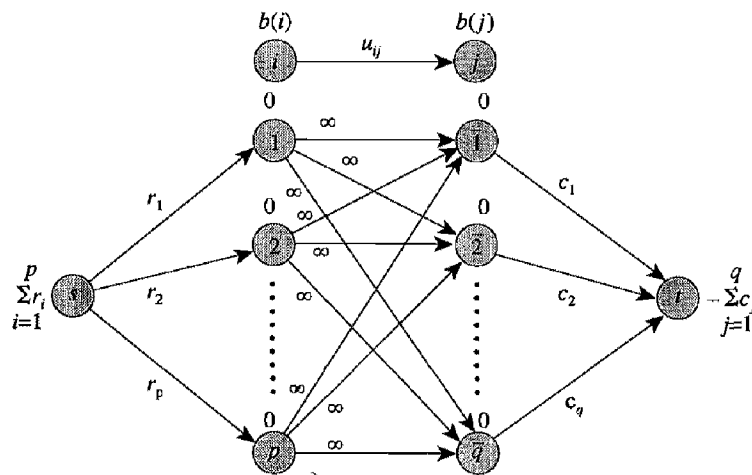


Figure 8.18 Feasible flow network for ensuring the statistical security of data.

send $k - x_{ij}^*$ units of flow from node j to node i and increase the flow in arc (i, j) by $k - x_{ij}^*$, leading to a total of k units.

Thus, to find β_{ij} , we need to compute the maximum flow from node j to node i in $G(x^*)$ for each arc $(i, j) \in A$. Equivalently, we want to find the minimum capacity cut $f[j, i]$ from node j to node i in $G(x^*)$. As stated previously, we cannot find $f[i, j]$ for all node pairs $[i, j]$ in a directed network; but we can determine $\min\{f[i, j], f[j, i]\}$. We now use the fact that we need to compute $f[j, i]$ when $(i, j) \in A$, in which case $f[i, j] = \infty$ (because the network is uncapacitated). Therefore, for each arc $(i, j) \in A$, $f[i, j] = \min\{f[i, j], f[j, i]\}$, and we can use the all-pairs minimum value cut algorithm to determine all of the β_{ij} values by solving $(2n - 2)$ maximum flow problems.

8.8 SUMMARY

As we have noted in our study of shortest path problems in our earlier chapters, we can sometimes develop more efficient algorithms by restricting the class of networks that we wish to consider (e.g., networks with nonnegative costs, or acyclic networks). In this chapter we have developed efficient special purpose algorithms for several maximum flow problems with specialized structure: (1) unit capacity networks, in which all arcs have unit capacities; (2) unit capacity simple networks, in which each node has a single incoming or outgoing arc; (3) bipartite networks; and (4) planar networks. We also considered one specialization of the maximum flow problem, the problem of finding the maximum number of arc-disjoint paths between two nodes in a network, and one generalization of the minimum cut problem, finding minimum value cuts between all pairs of nodes. For this last problem we permitted ourselves to measure the value of any cut by a function that is more general than the sum of the capacities of the arcs in the cut. Finally, we considered one other advanced topic: the use of the dynamic trees data structure to efficiently implement the shortest augmenting path algorithm for the maximum flow problem.

Figure 8.19, which summarizes the basic features of these various algorithms,

Algorithm	Running time	Features
Maximum flow algorithm for unit capacity networks	$O(\min\{n^{2/3}m, m^{3/2}\})$	<ol style="list-style-type: none"> 1. Fastest available algorithm for solving the maximum flow problem in unit capacity networks. 2. Uses two phases: first applies the shortest augmenting path algorithm until the distance label of node s satisfies the condition $d(s) \geq d^* = \min\{\lceil 2n^{2/3} \rceil, \lceil m^{1/2} \rceil\}$. At this point, uses the labeling algorithm until it establishes a maximum flow. 3. Easy to implement and is likely to be efficient in practice.
Maximum flow algorithm for unit capacity simple networks	$O(n^{1/2}m)$	<ol style="list-style-type: none"> 1. Fastest available algorithm for solving the maximum flow problem in unit capacity simple networks. 2. Same two phase approach as the preceding algorithm, except $d^* = \lceil n^{1/2} \rceil$.

Figure 8.19 Summary of algorithms discussed in this chapter.

Algorithm	Running time	Features
Bipartite preflow-push algorithm	$O(n^2m)$	<ol style="list-style-type: none"> 1. Faster approach for solving maximum flow problems in bipartite networks satisfying the condition $n_1 < n_2$. 2. Improved implementation of the generic preflow-push algorithm discussed in Section 7.6. 3. Uses "two-arc" push rule in which we always push flow from an active node over two consecutive admissible arcs. 4. As discussed in the exercises, significant further improvements are possible if we examine active nodes in some specific order.
Planar maximum flow algorithm	$O(n \log n)$	<ol style="list-style-type: none"> 1. Highly efficient algorithm for solving the maximum flow problem in $s-t$ planar networks. 2. Constructs the dual network and solves a shortest path problem over it. The shortest path in the dual network yields a minimum cut in the original network and the shortest path distances yield a maximum flow. 3. Applicable only to undirected $s-t$ planar networks.
Dynamic tree algorithm	$O(nm \log n)$	<ol style="list-style-type: none"> 1. Uses the dynamic tree data structure to implement the shortest augmenting path algorithm for the maximum flow problem. 2. Improves the running time of the shortest augmenting path algorithm from $O(n^2m)$ to $O(nm \log n)$. 3. Similar, though not as dramatic, improvements can be obtained by using this data structure in preflow-push algorithms. 4. The dynamic tree data structure is quite sophisticated, has substantial overhead and its practical usefulness has not yet been established.
Arc connectivity algorithm	$O(nm)$	<ol style="list-style-type: none"> 1. Fastest available algorithm for obtaining the arc connectivity of a network. 2. Uses the labeling algorithm for the maximum flow problem as a subroutine. 3. Likely to be very efficient in practice. 4. Applicable only to undirected networks.
All-pairs minimum cut algorithm	$O(nM(n, m, U) + n^3)$	<ol style="list-style-type: none"> 1. Fastest available algorithm for solving the all-pairs minimum cut problem. ($M(n, m, U)$ is the time needed for solving the maximum flow problem on a network with n nodes, m arcs, and U as the largest arc capacity.) 2. Determines minimum cuts between all pairs of nodes in the network by solving $(n - 1)$ maximum flow problems. 3. Can be used to determine the minimum value cuts between all pairs of nodes in the case in which we define the value of a cut differently than the capacity of the cut. 4. The $O(n^3)$ term in the worst-case bound can be reduced to $O(n^2)$ using different data structures. 5. Applicable to undirected networks only.

Figure 8.19 (Continued)

shows that by exploiting specialized structures or advanced data structures, we can improve on the running time of maximum flow computations, sometimes dramatically.

REFERENCE NOTES

We present the reference notes in this chapter separately for each of the several topics related to maximum flows that we have studied in this chapter.

Flows in unit capacity networks. Even and Tarjan [1975] showed that Dinic's algorithm solves the maximum flow problem in unit capacity and unit capacity simple networks in $O(\min\{n^{2/3}m, m^{3/2}\})$ and $O(n^{1/2}m)$ time, respectively. The algorithms we presented in Section 8.2 are due to Ahuja and Orlin [1991]; they use similar ideas and have the same running times. Fernandez-Baca and Martel [1989] presented and analyzed algorithms for solving more general maximum flow problems with "small" integer capacities.

Flows in bipartite networks. By improving on the running times of Dinic's [1970] and Karzanov's [1974] algorithms, Gusfield, Martel, and Fernandez-Baca [1987] developed the first specializations of maximum flow algorithms for bipartite networks. Ahuja, Orlin, Stein, and Tarjan [1990] provided further improvements and showed that it is possible to substitute n_1 for n in the time bounds of almost all preflow-push algorithms to obtain new time bounds for bipartite networks (recall that n_1 is the number of nodes on the smaller side of the bipartite network). This result implies that the generic preflow-push algorithm, the FIFO implementation, the highest-label implementation, and the excess scaling algorithm can solve the maximum flow problem in bipartite networks in $O(n_1^2m)$, $O(n_1m + n_1^3)$, $O(n_1m + n_1^2\sqrt{m})$, and $O(n_1m + n_1^2 \log U)$ time. Our discussion of the bipartite preflow-push algorithm in Section 8.3 is adapted from this paper. We have adapted the baseball elimination application from Schwartz [1966], and the network reliability application from Van Slyke and Frank [1972]. The paper by Gusfield, Martel, and Fernandez-Baca [1987] describes additional applications of bipartite maximum flow problems.

Flows in planar networks. In Section 8.4 we discussed the relationship between minimum s - t cuts in a network and shortest paths in its dual. Given a planar network G , the algorithm of Hopcroft and Tarjan [1974] constructs a planar representation in $O(n)$ time; from this representation, we can construct the dual network in $O(n)$ time. Berge [1957] showed that augmenting flow along certain paths, called *superior paths*, provides an algorithm that finds a maximum flow within n augmentations. Itai and Shiloach [1979] described an $O(n \log n)$ implementation of this algorithm. Hassin [1981] showed how to compute a maximum flow from the shortest path distances in the dual network. We have presented this method in our discussion in Section 8.4. For faster maximum flow algorithms in planar (but not necessarily s - t planar) undirected and directed networks, see Johnson and Venkatesan [1982] and Hassin and Johnson [1985].

Dynamic tree implementation. Sleator and Tarjan [1983] developed the dynamic tree data structure and used it to improve the worst-case complexity of Dinic's algorithm from $O(n^2m)$ to $O(nm \log n)$. Since then, researchers have used this data structure on many occasions to improve the performance of a range of network flow algorithms. Using the dynamic tree data structure, Goldberg and Tarjan [1986] improved the complexity of the FIFO preflow-push algorithm (described in Section 7.7) from $O(n^3)$ to $O(nm \log(n^2/m))$, and Ahuja, Orlin, and Tarjan [1989] improved the complexity of the excess scaling algorithm (described in Section 7.9) and several of its variants.

Network connectivity. Even and Tarjan [1975] offered an early discussion of arc connectivity of networks. Some of our discussion in Section 8.6 uses their results. The book by Even [1979] also contains a good discussion on node connectivity of a network. The $O(nm)$ time arc connectivity algorithm (for undirected networks) that we presented in Section 8.6 is due to Matula [1987] and is currently the fastest available algorithm. Mansour and Schieber [1988] presented an $O(nm)$ algorithm for determining the arc connectivity of a directed network.

All-pairs minimum value cut problem. Gomory and Hu [1961] developed the first algorithm for solving the all-pairs minimum cut problem on undirected networks that solves a sequence of $(n - 1)$ maximum flow problems. Gusfield [1990] presented an alternate all-pairs minimum cut algorithm that is very easy to implement using a code for the maximum flow problem. Talluri [1991] described yet a third approach. The algorithm we described in Section 8.7, which is due to Cheng and Hu [1990], is more general since it can handle cases when the value of a cut is defined differently than its capacity. Unfortunately, no one yet knows how to solve the all-pairs minimum value cut problem in directed networks as efficiently. No available algorithm is more efficient than solving $\Omega(n^2)$ maximum flow problems. The application of the all-pairs minimum value cut problem that we described at the end of Section 8.7 is due to Gusfield [1988]. Hu [1974] describes an additional application of the all-pairs minimum value cut problem that arises in network design.

EXERCISES

- 8.1 (a) Show that it is always possible to decompose a circulation in a unit capacity network into unit flows along arc-disjoint directed cycles.
- (b) Show that it is always possible to decompose a circulation in a simple network into unit flows along node-disjoint directed cycles.
- 8.2. Let $G = (N, A)$ be a directed network. Show that it is possible to decompose the arc set A into an arc-disjoint union of directed cycles if and only if G has a circulation x with $x_{ij} = 1$ for every arc $(i, j) \in A$. Moreover, show that we can find such a solution if and only if the indegree of each node equals its outdegree.
- 8.3. An undirected network is *biconnected* if it contains two node disjoint paths between every pair of nodes (except, of course, at the starting and terminal points). Show that a biconnected network must satisfy the following three properties: (1) for every two nodes p and q , and any arc (k, l) , some path from p to q contains arc (k, l) ; (2) for every three nodes p , q , and r , some path from p to r contains node q ; (3) for every three nodes p , q , and r , some path from p to r does not contain q .
- 8.4. Suppose that you are given a maximum flow problem in which all arc capacities are

the same. What is the most efficient method (from the worst-case complexity point of view) to solve this problem?

- 8.5. Using the unit capacity maximum flow algorithm, establish a maximum flow in the network shown in Figure 8.20.

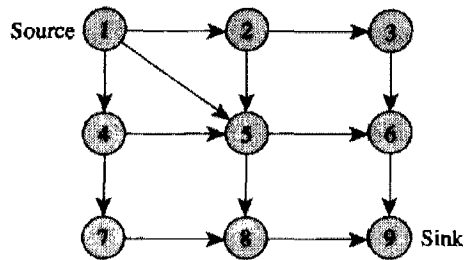


Figure 8.20 Example for Exercise 8.5.

- 8.6. Adapt the unit capacity maximum flow algorithm for unit capacity simple bipartite networks. In doing so, try to obtain the best possible running time. Describe your algorithm and analyze its worst-case complexity.
- 8.7. Consider a generalization of unit capacity networks in which arcs incident to the source and the sink nodes can have arbitrary capacities, but the remainder of the arcs have unit capacities. Will the unit capacity maximum flow algorithm still solve the problem in $O(\min\{n^{2/3}m, m^{3/2}\})$ time, or might the algorithm require more time? Consider a further generalization of the problem in which arcs incident to the source, the sink, and one other node have arbitrary capacities. What will be the complexity of the unit capacity maximum flow algorithm when applied to this problem?
- 8.8. We define a class of networks to be *small-capacity networks* if each arc capacity is between 1 and 4. Describe a generalization of the unit capacity maximum flow algorithm that would solve the maximum flow problems on small-capacity networks in $O(\min\{n^{2/3}m, m^{3/2}\})$ time.
- 8.9. What is the best possible bound you can obtain on the running time of the generic preflow-push algorithm applied to unit capacity networks?
- 8.10. Suppose we apply the preflow-push algorithm on a unit capacity simple network with the modification that we do not perform push/relabel operations on any node whose distance label exceeds $n^{1/2}$.
- Show that the modified preflow-push algorithm terminates within $O(n^{1/2}m)$ time.
 - Show that at the termination of the algorithm, the maximum additional flow that can reach the sink is at most $n^{1/2}$.
 - Can you convert this preflow into a maximum flow in $O(n^{1/2}m)$ time? If yes, then how?
- 8.11. Let x be a flow in a directed network. Assume that x is not a maximum flow. Let P and P' denote two successive shortest paths (i.e., P' is the shortest path after augmentation on path P) and suppose that P' contains at least one arc whose reversal lies in P . Show that $|P'| \geq |P| + 2$.
- 8.12. This exercise concerns the baseball elimination problem discussed in Application 8.1. Show that we can eliminate team 0 if and only if some nonempty set S of nodes satisfies the condition that

$$w_{\max} < \frac{\sum_{i \in S} w_i + \sum_{i \leq i < j \leq n} g_{ij} - \sum_{i \notin S \text{ and } j \in S} g_{ij}}{|S|}.$$

(Hint: Use the max-flow min-cut theorem.)

- 8.13. Given two n -element arrays α and β , we want to know whether we can construct an n -node directed graph so that node i has outdegree equal to $\alpha(i)$ and an indegree equal

to $\beta(i)$. Show how to solve this problem by solving a maximum flow problem. [Hint: Transform this problem to a feasible flow problem, as described in Application 6.1, on a complete bipartite graph $G = (N_1 \cup N_2, A)$ with $N_1 = \{1, 2, \dots, n\}$, $N_2 = \{1', 2', \dots, n'\}$, and $A = N_1 \times N_2$.]

- 8.14. Given an n -node graph G and two n -element arrays α and β , we wish to determine whether some subgraph G' of G satisfies the property that for each node i , $\alpha(i)$ and $\beta(i)$ are the outdegree and indegree of node i . Formulate this problem as a maximum flow problem. (Hint: The transformation is similar to the one used in Exercise 8.13.)
- 8.15. Apply the bipartite preflow-push algorithm to the maximum flow problem given in Figure 8.21. Among all the active nodes, push flow from a node with the smallest distance label.

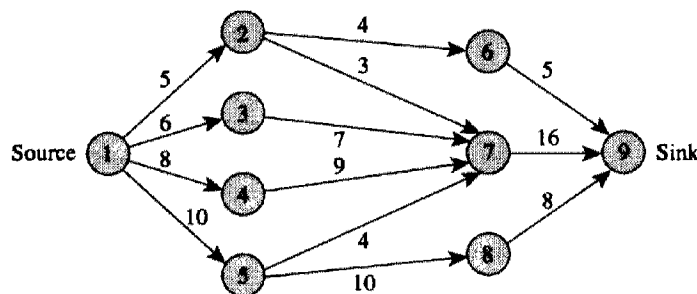


Figure 8.21 Example for Exercise 8.15.

- 8.16. Consider a bipartite network $G = (N_1 \cup N_2, A)$ with $n_1 = |N_1| \leq |N_2| = n_2$. Show that when applied to this network, the shortest augmenting path algorithm performs $O(n_1(n_1 + n_2)) = O(n_1 n_2)$ relabel steps and $O(n_1 m)$ augmentations, and runs in $O(n_1^2 m)$ time.
- 8.17. Suppose that we wish to find the maximum flow between two nodes in a bipartite network $G = (N_1 \cup N_2, A)$ with $n_1 = |N_1|$ and $n_2 = |N_2|$. This exercise considers the development of faster special implementations of the generic bipartite preflow-push algorithms.
- Show that if the algorithm always pushes flow from an active node with the highest distance label, it runs in $O(n_1^3 + n_1 m)$ time.
 - Show that if the algorithm examines active nodes in a FIFO order, it runs in $O(n_1^3 + n_1 m)$ time.
 - Develop an excess scaling version of the generic bipartite flow algorithm and show that it runs in $O(n_1 m + n_1^2 \log U)$ time.
- 8.18. A *semi-bipartite network* is defined as a network $G = (N, A)$ whose node set N can be partitioned into two subsets N_1 and N_2 so that no arc has both of its endpoints in N_2 (i.e., we allow arcs with both of their endpoints in N_1). Let $n_1 = |N_1|$, $n_2 = |N_2|$, and $n_1 \leq n_2$. Show how to modify the generic bipartite preflow-push algorithm so that it solves the maximum flow problem on semi-bipartite networks in $O(n_1^2 m)$ time.
- 8.19. (a) Prove that the graph shown in Figure 8.6(a) cannot be planar. (Hint: Use Euler's formula.)
 (b) Prove that the graph shown in Figure 8.6(b) cannot be planar. (Hint: Use Euler's formula.)
- 8.20. Show that an undirected planar network always contains a node with degree at most 5.
- 8.21. Apply the planar maximum flow algorithm to identify a maximum flow in the network shown in Figure 8.22.

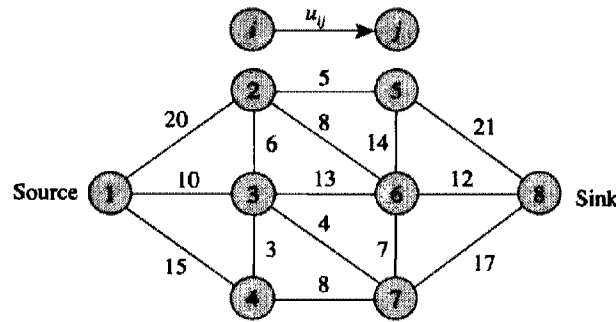


Figure 8.22 Example for Exercise 8.21.

- 8.22. Duals of directed s - t planar networks.** We define the dual graph of a directed s - t planar network $G = (N, A)$ as follows. We first draw an arc (t, s) of zero capacity, which divides the unbounded face into two faces: a new unbounded face and a new bounded face. We then place a node f^* inside each face f of the primal network G . Let s^* and t^* , respectively, denote the nodes in the dual network corresponding to the new bounded face and the new unbounded face. Each arc $(i, j) \in A$ lies on the boundary of the two faces f_1 and f_2 ; corresponding to this arc, the dual graph contains two oppositely directed arcs (f_1, f_2) and (f_2, f_1) . If arc (i, j) is a clockwise arc in the face f_1 , we define the cost of arc (f_1, f_2) as u_{ij} and the cost of (f_2, f_1) as zero. We define arc costs in the opposite manner if arc (i, j) is a counterclockwise arc in the face f_1 . Construct the dual of the s - t planar network shown in Figure 8.20. Next show that there is a one-to-one correspondence between s - t cuts in the primal network and directed paths from node s^* to node t^* in the dual network; moreover, show that the capacity of the cut equals the cost of the corresponding path.
- 8.23.** Show that if G is an s - t planar directed network, the minimum number of arcs in a directed path from s to t is equal to the maximum number of arc-disjoint s - t cuts. (*Hint:* Apply Theorem 6.7 to the dual of G .)
- 8.24. Node coloring algorithm.** In the node coloring problem, we wish to color the nodes of a network so that the endpoints of each arc have a different color. In this exercise we discuss an algorithm for coloring a planar undirected graph using at most six colors. The algorithm first orders the nodes of the network using the following iterative loop: It selects a node with degree at most 5 (from Exercise 8.20, we can always find such a node), deletes this node and its incident arcs from the network, and updates the degrees of all the nodes affected. The algorithm then examines nodes in the reverse order and assigns colors to them.
- Explain how to assign colors to the nodes to create a valid 6-coloring (i.e., the endpoints of every arc have a different color). Justify your method.
 - Show how to implement the node coloring algorithm so that it runs in $O(m)$ time.
- 8.25.** Consider the numerical example used in Section 8.5 to illustrate the dynamic tree implementation of the shortest augmenting path algorithm. Perform further iterations of the algorithm starting from Figure 8.11(c) until you find a maximum flow that has a value of 3.
- 8.26.** Solve the maximum flow problem given in Figure 8.23 by the dynamic tree implementation of the shortest augmenting path algorithm.

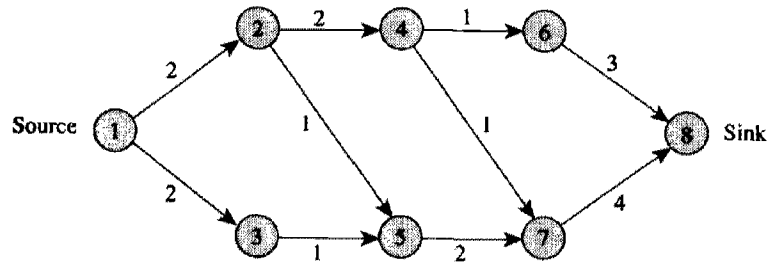


Figure 8.23 Example for Exercise 8.26.

- 8.27. Show how to use the dynamic tree data structure to implement in $O(m \log n)$ time the algorithm described in Exercise 7.11 for converting a maximum preflow into a maximum flow.
- 8.28. In Section 3.5 we showed how we can determine the flow decomposition of any flow in $O(nm)$ time. Show how to use the dynamic tree data structure to determine the flow decomposition in $O(m \log n)$ time.
- 8.29. Prove Property 8.11.
- 8.30. Compute the arc connectivity for the networks shown in Figure 8.24. Feel free to determine the maximum number of arc-disjoint paths between any pairs of nodes by inspection.

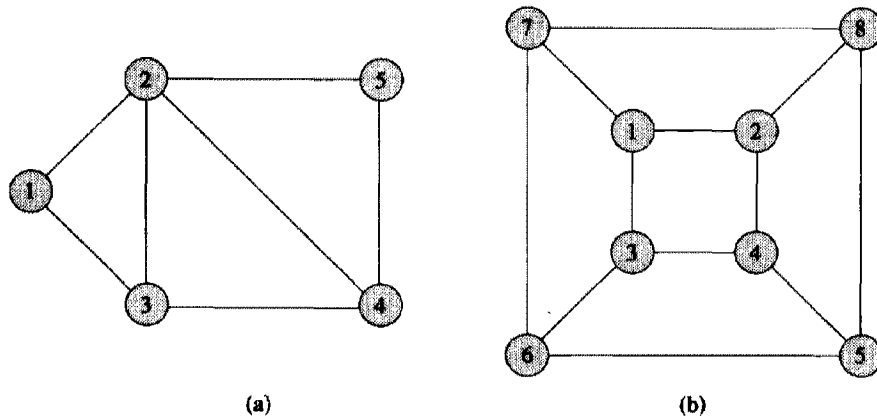


Figure 8.24 Example for Exercises 8.30 and 8.36.

- 8.31. Construct an undirected network whose nodes all have degree at least 3, but whose arc connectivity is 2.
- 8.32. An undirected network is said to be k -connected if every pair of nodes are connected by at least k arc-disjoint paths. Describe an $O(m)$ algorithm for determining whether a network is 1-connected. Use this algorithm to describe a simple $O(knm)$ algorithm for determining whether a network is k -connected. This algorithm should be different than those described in Section 8.6.
- 8.33. In a directed graph G , we define the *arc connectivity*, $\alpha[i, j]$, of an ordered pair $[i, j]$ of nodes as the minimum number of arcs whose deletion from the network eliminates all the directed path from node i to node j . We define the arc connectivity of a network G as $\alpha(G) = \min\{\alpha[i, j] : [i, j] \in N \times N\}$. Describe an algorithm for determining $\alpha(G)$ in $O(\min\{n^{5/3}m, m^2\})$ time and prove that this algorithm correctly determines the arc connectivity of any directed network. (Hint: Let p be a node in G of minimum degree. Determine $\alpha[p, j]$ and $\alpha[j, p]$ for all j , and take the minimum of these numbers.)

- 8.34. Arc connectivity of directed networks** (Schnorr [1979]). In Exercise 8.33 we showed how to determine the arc connectivity $\alpha(G)$ of a directed network G by solving at most $2n$ maximum flow problems. Prove the following result, which would enable us to determine the arc connectivity by solving n maximum flow problems. Let $1, 2, \dots, n$ be any ordering of the nodes in the network, and let node $(n + 1) = 1$. Show that $\alpha(G) = \min\{\alpha[i, i + 1] : i = 1, \dots, n\}$.
- 8.35. Node connectivity of undirected networks.** We define the *node connectivity*, $\beta[i, j]$, of a pair $[i, j]$ of nodes in an undirected graph $G = (N, A)$ as the minimum number of nodes whose deletion from the network eliminates all directed paths from node i to node j .
- Show that if $(i, j) \in A$, then $\beta[i, j]$ is not defined.
 - Let $H = \{(i, j) \in N \times N : (i, j) \notin A\}$ and let $\beta(G) = \min\{\beta[i, j] : [i, j] \in H\}$ denote the node connectivity of a network G . Show that $\beta(G) \leq 2 \lfloor m/n \rfloor$.
 - Show that the node connectivity of a network is no more than its arc connectivity.
 - A natural strategy for determining the node connectivity of a network would be to generalize an arc connectivity algorithm described in Section 8.6. We fix a node p and determine $\beta[p, j]$ for each j for which $(p, j) \notin A$. Using Figure 8.25 show that the minimum of these values will not be equal to $\beta(G)$. Explain why this approach fails for finding node connectivity even though it works for finding arc connectivity.

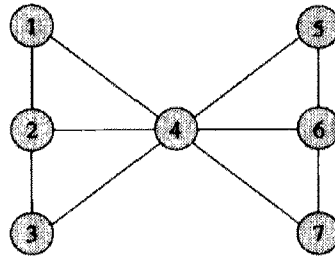


Figure 8.25 Example for Exercise 8.35.

- 8.36.** Solve the all-pairs minimum cut problems given in Figure 8.24. Obtain the separator tree for each problem.

9

MINIMUM COST FLOWS: BASIC ALGORITHMS

. . . men must walk, at least, before they dance.
—Alexander Pope

Chapter Outline

9.1	Introduction
9.2	Applications
9.3	Optimality Conditions
9.4	Minimum Cost Flow Duality
9.5	Relating Optimal Flows to Optimal Node Potentials
9.6	Cycle-Canceling Algorithm and the Integrality Property
9.7	Successive Shortest Path Algorithm
9.8	Primal–Dual Algorithm
9.9	Out-of-Kilter Algorithm
9.10	Relaxation Algorithm
9.11	Sensitivity Analysis
9.12	Summary

9.1 INTRODUCTION

The minimum cost flow problem is the central object of study in this book. In the last five chapters, we have considered two special cases of this problem: the shortest path problem and the maximum flow problem. Our discussion has been multifaceted: (1) We have seen how these problems arise in application settings as diverse as equipment replacement, project planning, production scheduling, census rounding, and analyzing round-robin tournaments; (2) we have developed a number of algorithmic approaches for solving these problems and studied their computational complexity; and (3) we have shown connections between these problems and more general problems in combinatorial optimization such as the minimum cut problem and a variety of min-max duality results. As we have seen, it is easy to understand the basic nature of shortest path and maximum flow problems and to develop core algorithms for solving them; nevertheless, designing and analyzing efficient algorithms is a very challenging task, requiring considerable ingenuity and considerable insight concerning both basic algorithmic strategies and their implementations.

As we begin to study more general minimum cost flow problems, we might ask ourselves a number of questions.

1. How much more difficult is it to solve the minimum cost flow problem than its shortest path and maximum flow specializations?
2. Can we use some of the same basic algorithmic strategies, such as label-setting and label-correcting methods, and the many variants of the augmenting path methods (e.g., shortest augmenting paths, scaling methods) for solving minimum cost flow problems?
3. The shortest path problem and the maximum flow problem address different components of the overall minimum cost flow problem: Shortest path problems consider arc flow costs but no flow capacities; maximum flow problems consider capacities but only the simplest cost structure. Since the minimum cost flow problem combines these problem ingredients, can we somehow combine the material that we have examined for shortest path and maximum flow problems to develop optimality conditions, algorithms, and underlying theory for the minimum cost flow problem?

In this and the next two chapters, we provide (partial) answers to these questions. We develop a number of algorithms for solving the minimum cost flow problem. Although these algorithms are not as efficient as those for the shortest path and maximum flow problems, they still are quite efficient, and indeed, are among the most efficient algorithms known in applied mathematics, computer science, and operations research for solving large-scale optimization problems.

We also show that we can develop considerable insight and useful tools and methods of analysis by drawing on the material that we have developed already. For example, in order to give us a firm foundation for developing algorithms for solving minimum cost flow problems, in Section 9.3 we establish optimality conditions for minimum cost flow problems based on the notion of node potentials associated with the nodes in the underlying network. These node potentials are generalizations of the concept of distance labels that we used in our study of shortest path problems. Recall that we were able to use distance labels to characterize optimal shortest paths; in addition, we used the distance label optimality conditions as a starting point for developing the basic iterative label-setting and label-correcting algorithms for solving shortest path problems. We use the node potential in a similar fashion for minimum cost flow problems. The connection with shortest paths is much deeper, however, than this simple analogy between node potentials and distance labels. For example, we show how to interpret and find the optimal node potentials for a minimum cost flow problem by solving an appropriate shortest path problem: The optimal node potentials are equal to the negative of the optimal distance labels from this shortest path problem.

In addition, many algorithms for solving the minimum cost flow problem combine ingredients of both shortest path and maximum flow algorithms. Many of these algorithms solve a sequence of shortest path problems with respect to maximum flow-like residual networks and augmenting paths. (Actually, to define the residual network, we consider both cost and capacity considerations.) We consider four such algorithms in this chapter. The *cycle-canceling algorithm* uses shortest path computations to find augmenting cycles with negative flow costs; it then augments flows along these cycles and iteratively repeats these computations for detecting negative

cost cycles and augmenting flows. The *successive shortest path algorithm* incrementally loads flow on the network from some source node to some sink node, each time selecting an appropriately defined shortest path. The *primal-dual* and *out-of-kilter algorithms* use a similar algorithmic strategy: at every iteration, they solve a shortest path problem and augment flow along one or more shortest paths. They vary, however, in their tactics. The primal-dual algorithm uses a maximum flow computation to augment flow simultaneously along several shortest paths. Unlike all the other algorithms, the out-of-kilter algorithm permits arc flows to violate their flow bounds. It uses shortest path computations to find flows that satisfy both the flow bounds and the cost and capacity based optimality conditions.

The fact that we can implement iterative shortest path algorithms in so many ways demonstrates the versatility that we have in solving minimum cost flow problems. Indeed, as we shall see in the next two chapters, we have even more versatility. Each of the algorithms that we discuss in this chapter is pseudopolynomial for problems with integer data. As we shall see in Chapter 10, by using ideas such as scaling of the problem data, we can also develop polynomial-time algorithms.

Since minimum cost flow problems are linear programs, it is not surprising to discover that we can also use linear programming methodologies to solve minimum cost flow problems. Indeed, many of the various optimality conditions that we have introduced in previous chapters and that we consider in this chapter are special cases of the more general optimality conditions of linear programming. Moreover, we can interpret many of these results in the context of a general theory of duality for linear programs. In this chapter we develop these duality results for minimum cost flow problems. In Chapter 11 we study the application of the key algorithmic approach from linear programming, the simplex method, for the minimum cost flow problem. In this chapter we consider one other algorithm, known as the *relaxation algorithm*, for solving the minimum cost flow problem.

To begin our discussion of the minimum cost flow problem, we first consider some additional applications, which help to show the importance of this problem in practice. Before doing so, however, let us set our notation and some underlying definitions that we use throughout our discussion.

Notation and Assumptions

Let $G = (N, A)$ be a directed network with a *cost* c_{ij} and a *capacity* u_{ij} associated with every arc $(i, j) \in A$. We associate with each node $i \in N$ a number $b(i)$ which indicates its supply or demand depending on whether $b(i) > 0$ or $b(i) < 0$. The minimum cost flow problem can be stated as follows:

$$\text{Minimize } z(x) = \sum_{(i,j) \in A} c_{ij}x_{ij} \quad (9.1a)$$

subject to

$$\sum_{(j:(i,j) \in A)} x_{ij} - \sum_{(j:(j,i) \in A)} x_{ji} = b(i) \quad \text{for all } i \in N, \quad (9.1b)$$

$$0 \leq x_{ij} \leq u_{ij} \quad \text{for all } (i, j) \in A. \quad (9.1c)$$

Let C denote the largest magnitude of any arc cost. Further, let U denote the

largest magnitude of any supply/demand or finite arc capacity. We assume that the lower bounds l_{ij} on arc flows are all zero. We further make the following assumptions:

Assumption 9.1. *All data (cost, supply/demand, and capacity) are integral.*

As noted previously, this assumption is not really restrictive in practice because computers work with rational numbers which we can convert to integer numbers by multiplying by a suitably large number.

Assumption 9.2. *The network is directed.*

We have shown in Section 2.4 that we can always fulfill this assumption by transforming any undirected network into a directed network.

Assumption 9.3. *The supplies/demands at the nodes satisfy the condition $\sum_{i \in N} b(i) = 0$ and the minimum cost flow problem has a feasible solution.*

We can determine whether the minimum cost flow problem has a feasible solution by solving a maximum flow problem as follows. Introduce a source node s^* and a sink node t^* . For each node i with $b(i) > 0$, add a “source” arc (s^*, i) with capacity $b(i)$, and for each node i with $b(i) < 0$, add a “sink” arc (i, t^*) with capacity $-b(i)$. Now solve a maximum flow problem from s^* to t^* . If the maximum flow saturates all the source arcs, the minimum cost flow problem is feasible; otherwise, it is infeasible. For the justification of this method, see Application 6.1 in Section 6.2.

Assumption 9.4. *We assume that the network G contains an uncapacitated directed path (i.e., each arc in the path has infinite capacity) between every pair of nodes.*

We impose this condition, if necessary, by adding *artificial* arcs $(1, j)$ and $(j, 1)$ for each $j \in N$ and assigning a large cost and infinite capacity to each of these arcs. No such arc would appear in a minimum cost solution unless the problem contains no feasible solution without artificial arcs.

Assumption 9.5. *All arc costs are nonnegative.*

This assumption imposes no loss of generality since the arc reversal transformation described in Section 2.4 converts a minimum cost flow problem with negative arc lengths to those with nonnegative arc lengths. This transformation, however, requires that all arcs have finite capacities. When some arcs are uncapacitated, we assume that the network contains no directed negative cost cycle of infinite capacity. If the network contains any such cycles, the optimal value of the minimum cost flow problem is unbounded; moreover, we can detect such a situation by using the search algorithm described in Section 3.4. In the absence of a negative cycle with infinite capacity, we can make each uncapacitated arc capacitated by setting its capacity equal to B , where B is the sum of all arc capacities and the supplies of all supply nodes; we justify this transformation in Exercise 9.36.

Residual Network

Our algorithms rely on the concept of residual networks. The residual network $G(x)$ corresponding to a flow x is defined as follows. We replace each arc $(i, j) \in A$ by two arcs (i, j) and (j, i) . The arc (i, j) has cost c_{ij} and *residual capacity* $r_{ij} = u_{ij} - x_{ij}$, and the arc (j, i) has cost $c_{ji} = -c_{ij}$ and residual capacity $r_{ji} = x_{ij}$. The residual network consists *only* of arcs with positive residual capacity.

9.2 APPLICATIONS

Minimum cost flow problems arise in almost all industries, including agriculture, communications, defense, education, energy, health care, manufacturing, medicine, retailing, and transportation. Indeed, minimum cost flow problems are pervasive in practice. In this section, by considering a few selected applications that arise in distribution systems planning, medical diagnosis, public policy, transportation, manufacturing, capacity planning, and human resource management, we give a passing glimpse of these applications. This discussion is intended merely to introduce several important applications and to illustrate some of the possible uses of minimum cost flow problems in practice. Taken together, the exercises in this chapter and in Chapter 11 and the problem descriptions in Chapter 19 give a much more complete picture of the full range of applications of minimum cost flows.

Application 9.1 Distribution Problems

A large class of network flow problems centers around shipping and distribution applications. One core model might be best described in terms of shipments from plants to warehouses (or, alternatively, from warehouses to retailers). Suppose that a firm has p plants with known supplies and q warehouses with known demands. It wishes to identify a flow that satisfies the demands at the warehouses from the available supplies at the plants and that minimizes its shipping costs. This problem is a well-known special case of the minimum cost flow problem, known as the *transportation problem*. We next describe in more detail a slight generalization of this model that also incorporates manufacturing costs at the plants.

A car manufacturer has several manufacturing plants and produces several car models at each plant that it then ships to geographically dispersed retail centers throughout the country. Each retail center requests a specific number of cars of each model. The firm must determine the production plan of each model at each plant and a shipping pattern that satisfies the demands of each retail center and minimizes the overall cost of production and transportation.

We describe this formulation through an example. Figure 9.1 illustrates a situation with two manufacturing plants, two retailers, and three car models. This model has four types of nodes: (1) *plant nodes*, representing various plants; (2) *plant/model nodes*, corresponding to each model made at a plant; (3) *retailer/model nodes*, corresponding to the models required by each retailer; and (4) *retailer nodes* corresponding to each retailer. The network contains three types of arcs.

1. *Production arcs*. These arcs connect a plant node to a plant/model node; the cost of this arc is the cost of producing the model at that plant. We might place

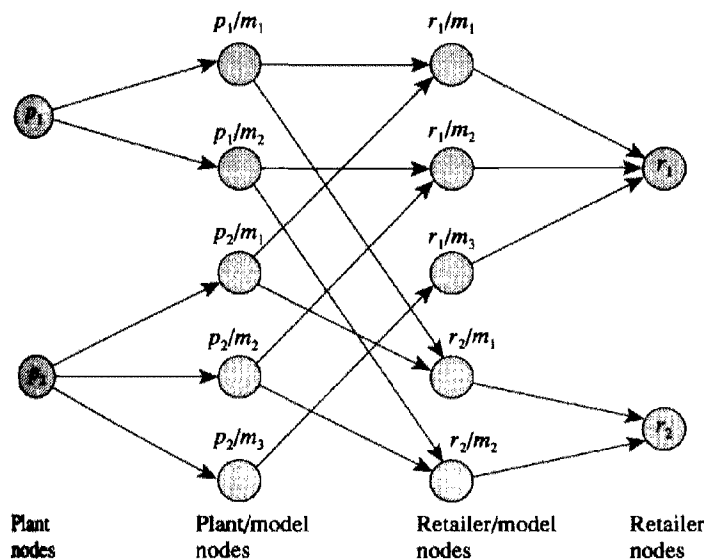


Figure 9.1 Production-distribution model.

lower and upper bounds on these arcs to control for the minimum and maximum production of each particular car model at the plants.

2. *Transportation arcs.* These arcs connect plant/model nodes to retailer/model nodes; the cost of such an arc is the total cost of shipping one car from the manufacturing plant to the retail center. Any such arc might correspond to a complex distribution channel with, for example, three legs: (a) a delivery from a plant (by truck) to a rail system; (b) a delivery from the rail station to another rail station elsewhere in the system; and (c) a delivery from the rail station to a retailer (by a local delivery truck). The transportation arcs might have lower or upper bounds imposed on their flows to model contractual agreements with shippers or capacities imposed on any distribution channel.
3. *Demand arcs.* These arcs connect retailer/model nodes to the retailer nodes. These arcs have zero costs and positive lower bounds which equal the demand of that model at that retail center.

Clearly, the production and shipping schedules for the automobile company correspond in a one-to-one fashion with the feasible flows in this network model. Consequently, a minimum cost flow would yield an optimal production and shipping schedule.

Application 9.2 Reconstructing the Left Ventricle from X-ray Projections

This application describes a network flow model for reconstructing the three-dimensional shape of the left ventricle from biplane angiocardiograms that the medical profession uses to diagnose heart diseases. To conduct this analysis, we first reduce the three-dimensional reconstruction problem into several two-dimensional problems by dividing the ventricle into a stack of parallel cross sections. Each two-dimensional cross section consists of one connected region of the left ventricle.

During a cardiac catheterization, doctors inject a dye known as Roentgen contrast agent into the ventricle; by taking x-rays of the dye, they would like to determine what portion of the left ventricle is functioning properly (i.e., permitting the flow of blood). Conventional biplane x-ray installations do not permit doctors to obtain a complete picture of the left ventricle; rather, these x-rays provide one-dimensional projections that record the total intensity of the dye along two axes (see Figure 9.2). The problem is to determine the distribution of the cloud of dye within the left ventricle and thus the shape of the functioning portion of the ventricle, assuming that the dye mixes completely with the blood and fills the portions that are functioning properly.

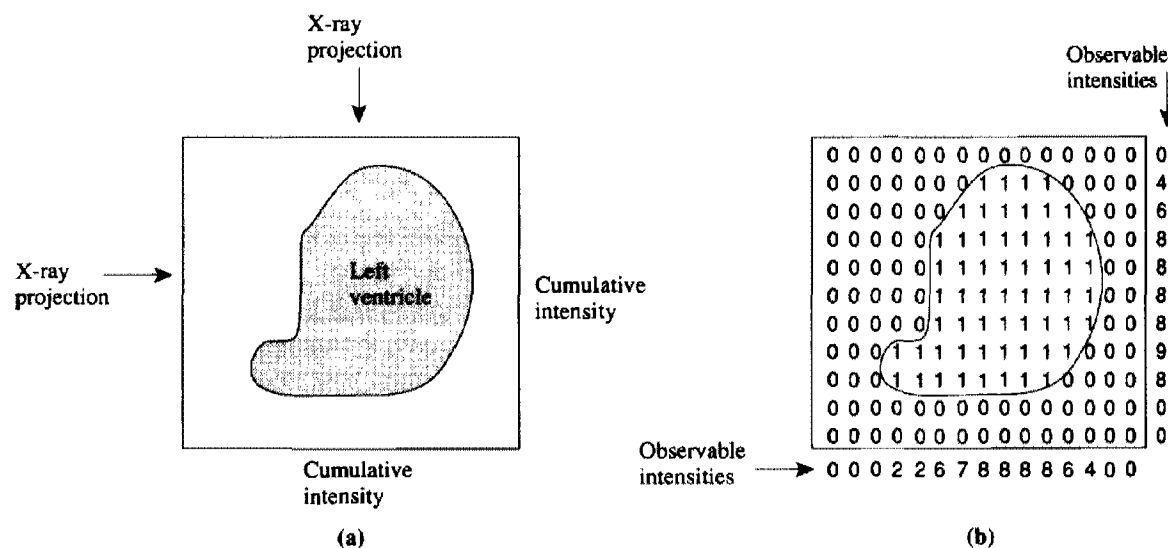


Figure 9.2 Using x-ray projections to measure a left ventricle.

We can conceive of a cross section of the ventricle as a $p \times r$ binary matrix: a 1 in a position indicates that the corresponding segment allows blood to flow and a 0 indicates that it does not permit blood to flow. The angiocardiograms give the cumulative intensity of the contrast agent in two planes which we can translate into row and column sums of the binary matrix. The problem is then to construct the binary matrix given its row and column sums. This problem is a special case of the feasible flow problem that we discussed in Section 6.2.

Typically, the number of feasible solutions for such problems are quite large; and these solutions might differ substantially. To constrain the feasible solutions, we might use certain facts from our experience that indicate that a solution is more likely to contain certain segments rather than others. Alternatively, we can use *a priori* information: for example, after some small time interval, the cross sections might resemble cross sections determined in a previous examination. Consequently, we might attach a probability p_{ij} that a solution will contain an element (i, j) of the binary matrix and might want to find a feasible solution with the largest possible cumulative probability. This problem is equivalent to a minimum cost flow problem.

Application 9.3 Racial Balancing of Schools

In Application 1.10 in Section 1.3 we formulated the racial balancing of schools as a multicommodity flow problem. We now consider a related, yet important situation: seeking a racial balance of two ethnic communities (blacks and whites). In this case we show how to formulate the problem as a minimum cost flow problem.

As in Application 1.10, suppose that a school district has S schools. For the purpose of this formulation, we divide the school district into L district locations and let b_i and w_i denote the number of black and white students at location i . These locations might, for example, be census tracts, bus stops, or city blocks. The only restrictions on the locations is that they be finite in number and that there be a single distance measure d_{ij} that reasonably approximates the distance any student at location i must travel if he or she is assigned to school j . We make the reasonable assumption that we can compute the distances d_{ij} before assigning students to schools. School j can enroll u_j students. Finally, let \underline{p} denote a lower bound and \bar{p} denote an upper bound on the percentage of black students assigned to each school (we choose these numbers so that school j has same percentage of blacks as does the school district). The objective is to assign students to schools in a manner that maintains the stated racial balance and minimizes the total distance traveled by the students.

We model this problem as a minimum cost flow problem. Figure 9.3 shows the minimum cost flow network for a three-location, two-school problem. Rather than describe the general model formally, we merely describe the model ingredients for this figure. In this formulation we model each location i as two nodes l_i' and l_i'' and each school j as two nodes s_j' and s_j'' . The decision variables for this problem are

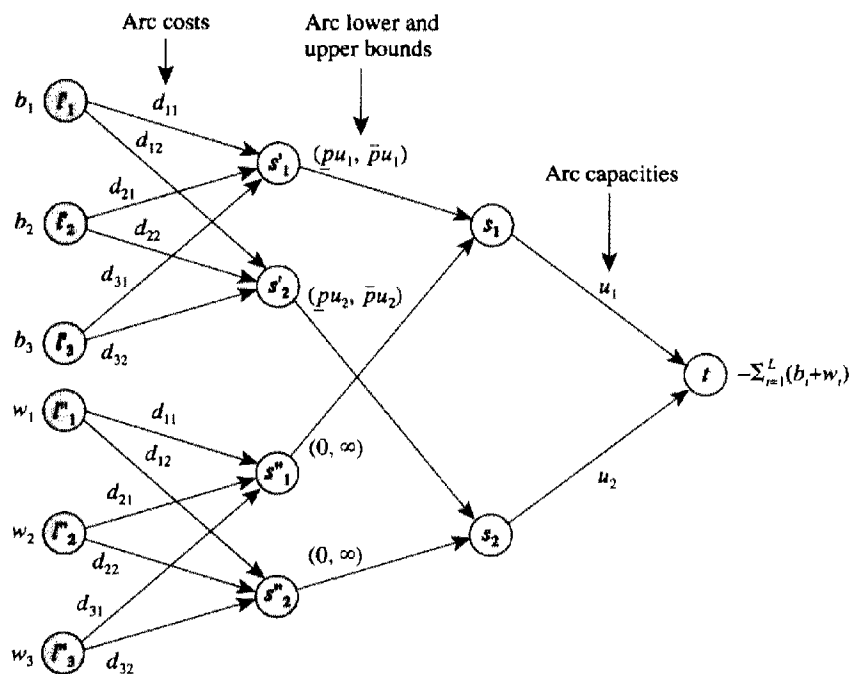


Figure 9.3 Network for the racial balancing of schools.

the number of black students assigned from location i to school j (which we represent by an arc from node l'_i to node s'_j) and the number of white students assigned from location i to school j (which we represent by an arc from node l''_i to node s''_j). These arcs are uncapacitated and we set their per unit flow cost equal to d_{ij} . For each j , we connect the nodes s'_j and s''_j to the school node s_j . The flow on the arcs (s'_j, s_j) and (s''_j, s_j) denotes the total number of black and white students assigned to school j . Since each school must satisfy lower and upper bounds on the number of black students it enrolls, we set the lower and upper bounds of the arc (s'_j, s_j) equal to $(\underline{p}u_j, \bar{p}u_j)$. Finally, we must satisfy the constraint that school j enrolls at most u_j students. We incorporate this constraint in the model by introducing a sink node t and joining each school node j to node t by an arc of capacity u_j . As is easy to verify, this minimum cost flow problem correctly models the racial balancing application.

Application 9.4 Optimal Loading of a Hopping Airplane

A small commuter airline uses a plane, with a capacity to carry at most p passengers, on a “hopping flight,” as shown in Figure 9.4(a). The hopping flight visits the cities $1, 2, 3, \dots, n$, in a fixed sequence. The plane can pick up passengers at any node and drop them off at any other node. Let b_{ij} denote the number of passengers available at node i who want to go to node j , and let f_{ij} denote the fare per passenger from node i to node j . The airline would like to determine the number of passengers that the plane should carry between the various origins to destinations in order to maximize the total fare per trip while never exceeding the plane capacity.

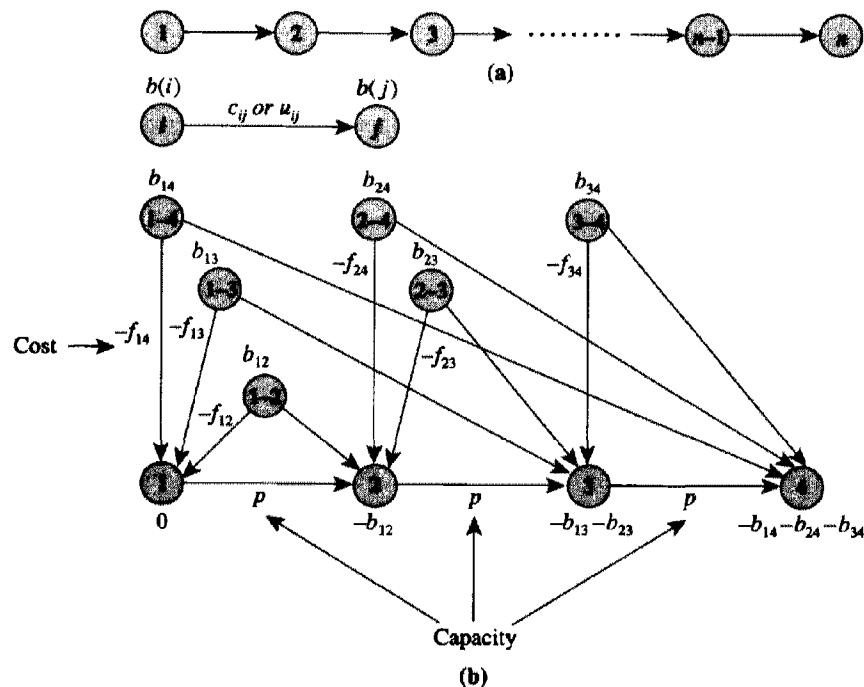


Figure 9.4 Formulating the hopping plane flight problem as a minimum cost flow problem.

Figure 9.4(b) shows a minimum cost flow formulation of this hopping plane flight problem. The network contains data for only those arcs with nonzero costs and with finite capacities: Any arc without an associated cost has a zero cost; any arc without an associated capacity has an infinite capacity. Consider, for example, node 1. Three types of passengers are available at node 1, those whose destination is node 2, node 3, or node 4. We represent these three types of passengers by the nodes 1-2, 1-3, and 1-4 with supplies b_{12} , b_{13} , and b_{14} . A passenger available at any such node, say 1-3, either boards the plane at its origin node by flowing through the arc (1-3, 1), and thus incurring a cost of $-f_{13}$ units, or never boards the plane which we represent by the flow through the arc (1-3, 3). In Exercise 9.13 we ask the reader to show that this formulation correctly models the hopping plane application.

Application 9.5 Scheduling with Deferral Costs

In some scheduling applications, jobs do not have any fixed completion times, but instead incur a deferral cost for delaying their completion. Some of these scheduling problems have the following characteristics: one of q identical processors (machines) needs to process each of p jobs. Each job j has a fixed processing time α_j that does not depend on which machine processes the job, or which jobs precede or follow the job. Job j also has a *deferral cost* $c_j(\tau)$, which we assume is a monotonically nondecreasing function of τ , the completion time of the job. Figure 9.5(a) illustrates one such deferral cost function. We wish to find a schedule for the jobs, with completion times denoted by $\tau_1, \tau_2, \dots, \tau_p$, that minimizes the total deferral cost $\sum_{j=1}^p c_j(\tau_j)$. This scheduling problem is difficult if the jobs have different processing times, but can be modeled as a minimum cost flow problem for situations with uniform processing times (i.e., $\alpha_j = \alpha$ for each $j = 1, \dots, p$).

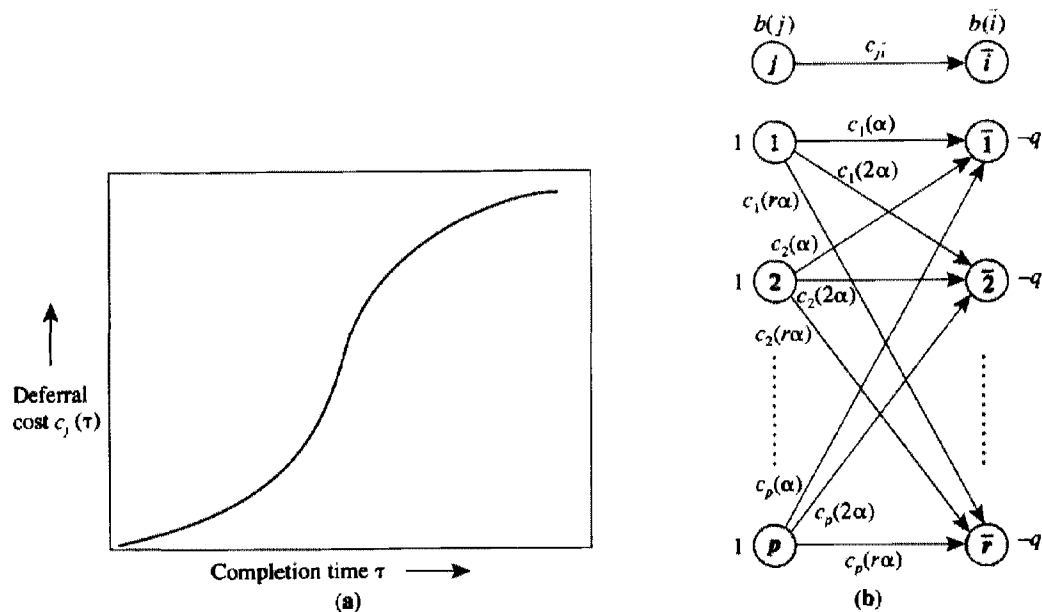


Figure 9.5 Formulating the scheduling problem with deferral costs.

Since the deferral costs are monotonically nondecreasing with time, in some optimal schedule the machines will process the jobs one immediately after another (i.e., the machines incur no *idle time*). As a consequence, in some optimal schedule the completion of each job will be $k\alpha$ for some constant k . The first job assigned to every machine will have a completion time of α units, the second job assigned to every machine will have a completion time of 2α units, and so on. This observation allows us to formulate the scheduling as a minimum cost flow problem in the network shown in Figure 9.5(b).

Assume, for simplicity, that $r = p/q$ is an integer. This assumption implies that we will assign exactly r jobs to each machine. (There is no loss of generality in imposing this assumption because we can add dummy jobs so that p/q becomes an integer.) The network has p job nodes, $1, 2, \dots, p$, each with 1 unit of supply; it also has r position nodes, $\bar{1}, \bar{2}, \dots, \bar{r}$, each with a demand of q units, indicating that the position has the capability to process q jobs. The flow on each arc (j, \bar{i}) is 1 or 0, depending on whether the schedule does or does not assign job j to the i th position of some machine. If we assign job j to the i th position on any machine, its completion time is $i\alpha$ and its deferral cost is $c_j(i\alpha)$. Therefore, arc (j, \bar{i}) has a cost of $c_j(i\alpha)$. Feasible schedules correspond, in a one-to-one fashion, with feasible flows in the network and both have the same cost. Consequently, a minimum cost flow will prescribe a schedule with the least possible deferral cost.

Application 9.6 Linear Programs with Consecutive 1's in Columns

Many linear programming problems of the form

$$\text{Minimize } cx$$

subject to

$$Ax \geq b,$$

$$x \geq 0,$$

have a special structure that permits us to solve the problem more efficiently than general-purpose linear programs. Suppose that the $p \times q$ matrix constraint matrix A is a 0–1 matrix satisfying the property that all of the 1's in each column appear consecutively (i.e., with no intervening zeros). We show how to transform this problem into a minimum cost flow problem. We illustrate our transformation using the following linear programming example:

$$\text{Minimize } cx \tag{9.2a}$$

subject to

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 \end{bmatrix} x \geq \begin{bmatrix} 5 \\ 12 \\ 10 \\ 6 \end{bmatrix}, \tag{9.2b}$$

$$x \geq 0. \tag{9.2c}$$

We first bring each constraint in (9.2b) into an equality form by introducing a “surplus” variable y_i for each row i in (9.2b). We then add a redundant row $0 \cdot x + 0 \cdot y = 0$ to the set of constraints. These changes produce the following equivalent formulation of the linear program:

$$\text{Minimize } cx \quad (9.3a)$$

subject to

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 1 & -1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 & 0 & -1 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & -1 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & -1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 5 \\ 12 \\ 10 \\ 6 \\ 0 \end{bmatrix}, \quad (9.3b)$$

$$x \geq 0. \quad (9.3c)$$

We next perform the following elementary row operation for each $i = p, p - 1, \dots, 1$, in the stated order: We subtract the i th constraint in (9.3b) from the $(i + 1)$ th constraint. These operations create the following equivalent linear program:

$$\text{Minimize } cx \quad (9.4a)$$

subject to

$$\begin{bmatrix} 0 & 1 & 0 & 1 & 1 & -1 & 0 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 & 1 & -1 & 0 & 0 \\ 0 & 0 & 1 & 0 & -1 & 0 & 1 & -1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & -1 \\ -1 & -1 & -1 & 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} 5 \\ 7 \\ -2 \\ -4 \\ -6 \end{bmatrix}, \quad (9.4b)$$

$$x \geq 0. \quad (9.4c)$$

Notice that in this form the constraints (9.4b) clearly define the mass balance constraints of a minimum cost flow problem because each column contains one $+1$ and one -1 . Also notice that the entries in the right-hand-side vector sum to zero, which is a necessary condition for feasibility. Figure 9.6 gives the minimum cost flow problem corresponding to this linear program.

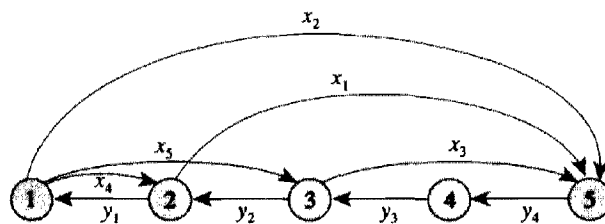


Figure 9.6 Formulating a linear program with consecutive ones as a minimum cost flow problem.

We have used a specific numerical example to illustrate the transformation of a linear program with consecutive 1's into a minimum cost flow problem. It is easy to show that this transformation is valid in general as well. For a linear program with p rows and q columns, the corresponding network has $p + 1$ nodes, one corresponding to each row, as well as one extra node that corresponds to an additional

“null row.” Each column A_k in the linear program that has consecutive 1’s in rows i to j becomes an arc $(i, j + 1)$ of cost c_k . Each surplus variable y_i becomes an arc $(i + 1, i)$ of zero cost. Finally, the supply/demand of a node i is $b(i) - b(i - 1)$.

Despite the fact that linear programs with consecutive 1’s might appear to be very special, and even contrived, this class of problems arises in a surprising number of applications. We illustrate the range of applications with three practical examples. We leave the formulations of these applications as minimum cost flow problems as exercises to the reader.

Optimal capacity scheduling. A vice-president of logistics of a large manufacturing firm must contract for $d(i)$ units of warehousing capacity for the time periods $i = 1, 2, \dots, n$. Let c_{ij} denote the cost of acquiring 1 unit of capacity at the beginning of period i , which is available for possible use throughout periods $i, i + 1, \dots, j - 1$ (assume that we relinquish this warehousing capacity at the beginning of period j). The vice-president wants to know how much capacity to acquire, at what times, and for how many subsequent periods, to meet the firm’s requirements at the lowest possible cost. This optimization problem arises because of possible savings that the firm might accrue by undertaking long-term leasing contracts at favorable times, even though these commitments might create excess capacity during some periods.

Employment scheduling. The vice-president of human resources of a large retail company must determine an employment policy that properly balances the cost of hiring, training, and releasing short-term employees, with the expense of having idle employees on the payroll for time periods when demand is low. Suppose that the company knows the minimum labor requirement d_j for each period $j = 1, \dots, n$. Let c_{ij} denote the cost of hiring someone at the beginning of period i and releasing him at the end of period $j - 1$. The vice-president would like to identify an employment policy that meets the labor requirements and minimizes the cost of hiring, training, and releasing employees.

Equipment replacement. A job shop must periodically replace its capital equipment because of machine wear. As a machine ages, it breaks down more frequently and so becomes more expensive to operate. Furthermore, as a machine ages, its salvage value decreases. Let c_{ij} denote the cost of buying a particularly important machine at the beginning of period i , plus the cost of operating the machine over the periods $i, i + 1, \dots, j - 1$, minus the salvage cost of the machine at the beginning of period j . The *equipment replacement problem* attempts to obtain a replacement plan that minimizes the total cost of buying, selling, and operating the machine over a planning horizon of n years, assuming that the job shop must have at least 1 unit of this machine in service at all times.

9.3 OPTIMALITY CONDITIONS

In our discussion of shortest path problems in Section 5.2, we saw that a set of distance labels $d(i)$ defines shortest path distances from a specified node s to every other node in the network if and only if they represent distances along some paths

from node s and satisfy the following shortest path optimality conditions:

$$d(j) \leq d(i) + c_{ij} \quad \text{for all } (i, j) \in A. \quad (9.5)$$

These optimality conditions are useful in several respects. First, they give us a simple validity check to see whether a given set of distance labels does indeed define shortest paths. Similarly, the optimality conditions provide us with a method for determining whether or not a given set of paths, one from node s to every other node in the network, constitutes a set of shortest paths from node s . We simply compute the lengths of these paths and see if these distances satisfy the optimality conditions. In both cases, the optimality conditions provide us with a “certificate” of optimality, that is, an assurance that a set of distance labels or a set of paths is optimal. One nice feature of the certificate is its ease of use. We need not invoke any complex algorithm to certify that a solution is optimal; we simply check the optimality conditions. The optimality conditions are also valuable for other reasons; as we saw in Chapter 5, they can suggest algorithms for solving a shortest path problem: For example, the generic label-correcting algorithm uses the simple idea of repeatedly replacing $d(j)$ by $d(i) + c_{ij}$ if $d(j) > d(i) + c_{ij}$ for some arc (i, j) . Finally, the optimality conditions provide us with a mechanism for establishing the validity of algorithms for the shortest path problem. To show that an algorithm correctly finds the desired shortest paths, we verify that the solutions they generate satisfy the optimality conditions.

These various uses of the shortest path optimality conditions suggest that similar sets of conditions might be valuable for designing and analyzing algorithms for the minimum cost flow problem. Accordingly, rather than launching immediately into a discussion of algorithms for solving the minimum cost flow problem, we first pause to describe a few different optimality conditions for this problem. All the optimality conditions that we state have an intuitive network interpretation and are rather direct extensions of their shortest path counterparts. We will consider three different (but equivalent) optimality conditions: (1) negative cycle optimality conditions, (2) reduced cost optimality conditions, and (3) complementary slackness optimality conditions.

Negative Cycle Optimality Conditions

The negative cycle optimality conditions stated next are a direct consequence of the flow decomposition property stated in Theorem 3.5 and our definition of residual networks given at the end of Section 9.1.

Theorem 9.1 (Negative Cycle Optimality Conditions). *A feasible solution x^* is an optimal solution of the minimum cost flow problem if and only if it satisfies the negative cycle optimality conditions: namely, the residual network $G(x^*)$ contains no negative cost (directed) cycle.*

Proof. Suppose that x is a feasible flow and that $G(x)$ contains a negative cycle. Then x cannot be an optimal flow, since by augmenting positive flow along the cycle we can improve the objective function value. Therefore, if x^* is an optimal flow, then $G(x^*)$ cannot contain a negative cycle. Now suppose that x^* is a feasible flow

and that $G(x^*)$ contains no negative cycle. Let x° be an optimal flow and $x^* \neq x^\circ$. The augmenting cycle property stated in Theorem 3.7 shows that we can decompose the difference vector $x^\circ - x^*$ into at most m augmenting cycles with respect to the flow x^* and the sum of the costs of flows on these cycles equals $cx^\circ - cx^*$. Since the lengths of all the cycles in $G(x^*)$ are nonnegative, $cx^\circ - cx^* \geq 0$, or $cx^\circ \geq cx^*$. Moreover, since x° is an optimal flow, $cx^\circ \leq cx^*$. Thus $cx^\circ = cx^*$, and x^* is also an optimal flow. This argument shows that if $G(x^*)$ contains no negative cycle, then x^* must be optimal, and this conclusion completes the proof of the theorem. ♦

Reduced Cost Optimality Conditions

To develop our second and third optimality conditions, let us make one observation. First, note that we can write the shortest path optimality conditions in the following equivalent form:

$$c_{ij}^d = c_{ij} + d(i) - d(j) \geq 0 \quad \text{for all arcs } (i, j) \in A. \quad (9.6)$$

This expression has the following interpretation: c_{ij}^d is an optimal “reduced cost” for arc (i, j) in the sense that it measures the cost of this arc relative to the shortest path distances $d(i)$ and $d(j)$. Notice that with respect to the optimal distances, every arc in the network has a nonnegative reduced cost. Moreover, since $d(j) = d(i) + c_{ij}$, if arc (i, j) is on a shortest path connecting the source node s to any other node, the shortest path uses only zero reduced cost arcs. Consequently, once we know the optimal distances, the problem is very easy to solve: We simply find a path from node s to every other node that uses only arcs with zero reduced costs. This interpretation raises a natural question: Is there a similar set of conditions for more general minimum cost flow problems?

Suppose that we associate a real number $\pi(i)$, unrestricted in sign, with each node $i \in N$. We refer to $\pi(i)$ as the *potential* of node i . We show in Section 9.4 that $\pi(i)$ is the linear programming dual variable corresponding to the mass balance constraint of node i . For a given set of node potentials π , we define the *reduced cost* of an arc (i, j) as $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$. These reduced costs are applicable to the residual network as well as the original network. We define the reduced costs in the residual network just as we did the costs, but now using c_{ij}^π in place of c_{ij} . The following properties will prove to be useful in our subsequent developments in this and later chapters.

Property 9.2

- (a) For any directed path P from node k to node l , $\sum_{(i,j) \in P} c_{ij}^\pi = \sum_{(i,j) \in P} c_{ij} - \pi(k) + \pi(l)$.
- (b) For any directed cycle W , $\sum_{(i,j) \in W} c_{ij}^\pi = \sum_{(i,j) \in W} c_{ij}$.

The proof of this property is similar to that of Property 2.5. Notice that this property implies that the node potentials do not change the shortest path between any pair of nodes k and l , since the potentials increase the length of every path by a constant amount $\pi(l) - \pi(k)$. This property also implies that if W is a negative cycle with respect to c_{ij} as arc costs, it is also a negative cycle with respect to c_{ij}^π .

as arc costs. We can now provide an alternative form of the negative cycle optimality conditions, stated in terms of the reduced costs of the arcs.

Theorem 9.3 (Reduced Cost Optimality Conditions). *A feasible solution x^* is an optimal solution of the minimum cost flow problem if and only if some set of node potentials π satisfy the following reduced cost optimality conditions:*

$$c_{ij}^\pi \geq 0 \quad \text{for every arc } (i, j) \text{ in } G(x^*). \quad (9.7)$$

Proof. We shall prove this result using Theorem 9.1. To show that the negative cycle optimality conditions is equivalent to the reduced cost optimality conditions, suppose that the solution x^* satisfies the latter conditions. Therefore, $\sum_{(i,j) \in W} c_{ij}^\pi \geq 0$ for every directed cycle W in $G(x^*)$. Consequently, by Property 9.2(b), $\sum_{(i,j) \in W} c_{ij}^\pi = \sum_{(i,j) \in W} c_{ij} \geq 0$, so $G(x^*)$ contains no negative cycle.

To show the converse, assume that for the solution x^* , $G(x^*)$ contains no negative cycle. Let $d(\cdot)$ denote the shortest path distances from node 1 to all other nodes in $G(x^*)$. Recall from Section 5.2 that if the network contains no negative cycle, the distance labels $d(\cdot)$ are well defined and satisfy the conditions $d(j) \leq d(i) + c_{ij}$ for all (i, j) in $G(x^*)$. We can restate these inequalities as $c_{ij} - (-d(i)) + (-d(j)) \geq 0$, or $c_{ij}^\pi \geq 0$ if we define $\pi = -d$. Consequently, the solution x^* satisfies the reduced cost optimality conditions. ♦

In the preceding theorem we characterized an optimal flow x as a flow that satisfied the conditions $c_{ij}^\pi \geq 0$ for all (i, j) in $G(x)$ for some set of node potentials π . In the same fashion, we could define “optimal node potentials” as a set of node potentials π that satisfy the conditions $c_{ij}^\pi \geq 0$ for all (i, j) in $G(x)$ for some feasible flow x .

We might note that the reduced cost optimality conditions have a convenient economic interpretation. Suppose that we interpret c_{ij} as the cost of transporting 1 unit of a commodity from node i to node j through the arc (i, j) , and we interpret $\mu(i) \equiv -\pi(i)$ as the cost of obtaining a unit of this commodity at node i . Then $c_{ij} + \mu(i)$ is the cost of the commodity at node j if we obtain it at node i and transport it to node j . The reduced cost optimality condition, $c_{ij} - \pi(i) + \pi(j) \geq 0$, or equivalently, $\mu(j) \leq c_{ij} + \mu(i)$, states that the cost of obtaining the commodity at node j is no more than the cost of the commodity if we obtain it at node i and incur the transportation cost in sending it from node i to j . The cost at node j might be smaller than $c_{ij} + \mu(i)$ because there might be a more cost-effective way to transport the commodity to node j via other nodes.

Complementary Slackness Optimality Conditions

Both Theorems 9.1 and 9.3 provide means for establishing optimality of solutions to the minimum cost flow problem by formulating conditions imposed on the residual network; we shall now restate these conditions in terms of the original network.

Theorem 9.4 (Complementary Slackness Optimality Conditions). *A feasible solution x^* is an optimal solution of the minimum cost flow problem if and only if for some set of node potentials π , the reduced costs and flow values satisfy the following complementary slackness optimality conditions for every arc $(i, j) \in A$:*

$$\text{If } c_{ij}^\pi > 0, \text{ then } x_{ij}^* = 0. \quad (9.8a)$$

$$\text{If } 0 < x_{ij}^* < u_{ij}, \text{ then } c_{ij}^\pi = 0. \quad (9.8b)$$

$$\text{If } c_{ij}^\pi < 0, \text{ then } x_{ij}^* = u_{ij}. \quad (9.8c)$$

Proof. We show that the reduced cost optimality conditions are equivalent to (9.8). To establish this result, we first prove that if the node potentials π and the flow vector x satisfy the reduced cost optimality conditions, then they must satisfy (9.8). Consider three possibilities for any arc $(i, j) \in A$.

Case 1. If $c_{ij}^\pi > 0$, the residual network cannot contain the arc (j, i) because $c_{ji}^\pi = -c_{ij}^\pi < 0$ for that arc, contradicting (9.7). Therefore, $x_{ij}^* = 0$.

Case 2. If $0 < x_{ij}^* < u_{ij}$, the residual network contains both the arcs (i, j) and (j, i) . The reduced cost optimality conditions imply that $c_{ij}^\pi \geq 0$ and $c_{ji}^\pi \geq 0$. But since $c_{ji}^\pi = -c_{ij}^\pi$, these inequalities imply that $c_{ij}^\pi = c_{ji}^\pi = 0$.

Case 3. If $c_{ij}^\pi < 0$, the residual network cannot contain the arc (i, j) because $c_{ij}^\pi < 0$ for that arc, contradicting (9.7). Therefore, $x_{ij}^* = u_{ij}$.

We have thus shown that if the node potentials π and the flow vector x satisfy the reduced cost optimality conditions, they also satisfy the complementary slackness optimality conditions. In Exercise 9.28 we ask the reader to prove the converse result: If the pair (x, π) satisfies the complementary slackness optimality conditions, it also satisfies the reduced cost optimality conditions. ♦

Those readers familiar with linear programming might notice that these conditions are the complementary slackness conditions for a linear programming problem whose variables have upper bounds; this association explains the choice of the name complementary slackness.

9.4 MINIMUM COST FLOW DUALITY

When we were introducing shortest path problems with nonnegative arc costs in Chapter 4, we considered a string model with knots representing the nodes of the network and with a string of length c_{ij} connecting the i th and j th knots. To solve the shortest path problem between a designated source node s and sink node t , we hold the string at the knots s and t and pull them as far apart as possible. As we noted in our previous discussion, if $d(i)$ denotes the distance from the source node s to node i along the shortest path and nodes i and j are any two nodes on this path, then $d(i) + c_{ij} \geq d(j)$. The shortest path distances might satisfy this inequality as a strict inequality if the string from node i to node j is not taut. In this string solution, since we are pulling the string apart as far as possible, we are obtaining the optimal shortest path distance between nodes s and t by solving a *maximization* problem. We could cast this problem formally as the following maximization problem:

$$\text{Maximize } d(t) - d(s) \quad (9.9a)$$

subject to

$$d(j) - d(i) \leq c_{ij} \quad \text{for all } (i, j) \in A. \quad (9.9b)$$

In this formulation, $d(s) = 0$. As we have noted in Chapter 4, if d is any vector of distance labels satisfying the constraints of this problem and the path P defined as $s = i_1 - i_2 - \cdots i_k - t$ is any path from node s to node t , then

$$\begin{aligned} d(i_1) - d(s) &\leq c_{si_1} \\ d(i_2) - d(i_1) &\leq c_{i_1 i_2} \\ &\vdots \\ d(t) - d(i_k) &\leq c_{i_k t}, \end{aligned}$$

so by adding these inequalities and using the fact that $d(s) = 0$, we see that

$$d(t) \leq c_{si_1} + c_{i_1 i_2} + \cdots + c_{i_k t}.$$

This result shows that if d is any feasible vector to the optimization problem (9.9), then $d(t)$ is a lower bound on the length of any path from node s to node t and therefore is a lower bound on the shortest distance between these nodes. As we see from the string solution, if we choose the distance labels $d(\cdot)$ appropriately (as the distances obtained from the string solution), $d(t)$ equals the shortest path distance.

This discussion shows the connection between the shortest path problem and a related maximization problem (9.9). In our discussion of the maximum flow problem, we saw a similar relationship, namely, the max-flow min-cut theorem, which tells us that associated with every maximum flow problem is an associated minimization problem. Moreover, since the maximum flow equals the minimum cut, the optimal value of these two associated problems is the same. These two results are special cases of a more general property that applies to any minimum cost flow problem, and that we now establish.

For every linear programming problem, which we subsequently refer to as a *primal* problem, we can associate another intimately related linear programming problem, called its *dual*. For example, the objective function value of any feasible solution of the dual is less than or equal to the objective function of any feasible solution of the primal. Furthermore, the maximum objective function value of the dual equals the minimum objective function of the primal. This duality theory is fundamental to an understanding of the theory of linear programming. In this section we state and prove these duality theory results for the minimum cost flow problem.

While forming the dual of a (primal) linear programming problem, we associate a *dual variable* with every constraint of the primal except for the nonnegativity restriction on arc flows. For the minimum cost flow problem stated in (9.1), we associate the variable $\pi(i)$ with the mass balance constraint of node i and the variable α_{ij} with the capacity constraint of arc (i, j) . In terms of these variables, the *dual minimum cost flow problem* can be stated as follows:

$$\text{Maximize } w(\pi, \alpha) = \sum_{i \in N} b(i)\pi(i) - \sum_{(i,j) \in A} u_{ij}\alpha_{ij} \quad (9.10a)$$

subject to

$$\pi(i) - \pi(j) - \alpha_{ij} \leq c_{ij} \quad \text{for all } (i, j) \in A, \quad (9.10b)$$

$$\alpha_{ij} \geq 0 \quad \text{for all } (i, j) \in A \quad \text{and} \quad \pi(j) \text{ unrestricted for all } j \in N. \quad (9.10c)$$

Note that the shortest path dual problem (9.9) is a special case of this model: For the shortest path problem, $b(s) = 1$, $b(t) = -1$, and $b(i) = 0$ otherwise. Also, since the shortest path problem contains no arc capacities, we can eliminate the α_{ij} variables. Therefore, if we let $d(i) = -\pi(i)$, the dual minimum cost flow problem (9.10) becomes the shortest path dual problem (9.9).

Our first duality result for the general minimum cost flow problem is known as the *weak duality theorem*.

Theorem 9.5 (Weak Duality Theorem). *Let $z(x)$ denote the objective function value of some feasible solution x of the minimum cost flow problem and let $w(\pi, \alpha)$ denote the objective function value of some feasible solution (π, α) of its dual. Then $w(\pi, \alpha) \leq z(x)$.*

Proof. We multiply both sides of (9.10b) by x_{ij} and sum these weighted inequalities for all $(i, j) \in A$, obtaining

$$\sum_{(i,j) \in A} (\pi(i) - \pi(j))x_{ij} - \sum_{(i,j) \in A} \alpha_{ij}x_{ij} \leq \sum_{(i,j) \in A} c_{ij}x_{ij}. \quad (9.11)$$

Notice that $cx - c^\pi x = \sum_{(i,j) \in A} (\pi(i) - \pi(j))x_{ij}$ [because $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$]. Next notice that Property 2.4 in Section 2.4 implies that $cx - c^\pi x$ equals $\sum_{i \in N} b(i)\pi(i)$. Therefore, the first term on the left-hand side of (9.11) equals $\sum_{i \in N} b(i)\pi(i)$. Next notice that replacing x_{ij} in the second term on the left-hand side of (9.11) by u_{ij} preserves the inequality because $x_{ij} \leq u_{ij}$ and $\alpha_{ij} \geq 0$. Consequently,

$$\sum_{i \in N} b(i)\pi(i) - \sum_{(i,j) \in A} \alpha_{ij}u_{ij} \leq \sum_{(i,j) \in A} c_{ij}x_{ij}. \quad (9.12)$$

Now notice that the left-hand side of (9.12) is the dual objective $w(\pi, \alpha)$ and the right-hand side is the primal objective, so we have established the lemma. ♦

The weak duality theorem implies that the objective function value of *any* dual feasible solution is a lower bound on the objective function value of *any* primal feasible solution. One consequence of this result is immediate: If some dual solution (π, α) and a primal solution x have the same objective function value, (π, α) must be an optimal solution of the dual problem and x must be an optimal solution of the primal problem (why?). Can we always find such solutions? The *strong duality theorem*, to be proved next, answers this question in the affirmative.

We first eliminate the dual variables α_{ij} 's from the dual formation (9.10) using some properties of the optimal solution. Defining the reduced cost, as before, as $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$, we can rewrite the constraint (9.10b) as

$$\alpha_{ij} \geq -c_{ij}^\pi. \quad (9.13)$$

The coefficient associated with the variable α_{ij} in the dual objective (9.10a) is $-u_{ij}$, and we wish to maximize the objective function value. Consequently, in any optimal solution we would assign the smallest possible value to α_{ij} . This observation, in view of (9.10c) and (9.13), implies that

$$\alpha_{ij} = \max\{0, -c_{ij}^\pi\}. \quad (9.14)$$

We have thus shown that if we know optimal values for the dual variables $\pi(i)$, we can compute the optimal values of the variables α_{ij} using (9.14). This construction permits us to eliminate the variables α_{ij} from the dual formulation. Substituting (9.14) in (9.10a) yields

$$\text{Maximize } w(\pi) = \sum_{i \in N} b(i)\pi(i) - \sum_{(i,j) \in A} \max\{0, -c_{ij}^{\pi}\}u_{ij}. \quad (9.15)$$

The dual problem reduces to finding a vector π that optimizes (9.15). We are now in a position to prove the strong duality theorem. (Recall that our blanket assumption, Assumption 9.3, implies that the minimum cost flow problem always has a solution.)

Theorem 9.6 (Strong Duality Theorem). *For any choice of problem data, the minimum cost flow problem always has a solution x^* and the dual minimum cost flow problem has a solution π satisfying the property that $z(x^*) = w(\pi)$.*

Proof. We prove this theorem using the complementary slackness optimality conditions (9.8). Let x^* be an optimal solution of the minimum cost flow problem. Theorem 9.4 implies that x^* together with some vector π of node potentials satisfy the complementary slackness optimality conditions. We claim that this solution satisfies the condition

$$-c_{ij}^{\pi}x_{ij}^* = \max\{0, -c_{ij}^{\pi}\}u_{ij} \quad \text{for every arc } (i, j) \in A. \quad (9.16)$$

To establish this result, consider the following three cases: (1) $c_{ij}^{\pi} > 0$, (2) $c_{ij}^{\pi} = 0$, and (3) $c_{ij}^{\pi} < 0$. The complementary slackness conditions (9.8) imply that in the first two cases, both the left-hand side and right-hand side of (9.16) are zero, and in the third case both sides equal $-c_{ij}^{\pi}u_{ij}$.

Next consider the dual objective (9.15). Substituting (9.16) in (9.15) yields

$$w(\pi) = \sum_{i \in N} b(i)\pi(i) + \sum_{(i,j) \in A} c_{ij}^{\pi}x_{ij}^* = \sum_{(i,j) \in A} c_{ij}x_{ij}^* = z(x^*).$$

The second last inequality follows from Property 2.4. This result is the conclusion of the theorem. ♦

The proof of this theorem shows that any optimal solution x^* of the minimum cost flow problem always has an associated dual solution π satisfying the condition $z(x^*) = w(\pi)$. Needless to say, the solution π is an optimal solution of the dual minimum cost flow problem since any larger value of the dual objective would contradict the weak duality theorem stated in Theorem 9.5.

In Theorem 9.6 we showed that the complementary slackness optimality conditions implies strong duality. We next prove the converse result: namely, that strong duality implies the complementary slackness optimality conditions.

Theorem 9.7. *If x is a feasible flow and π is an (arbitrary) vector satisfying the property that $z(x) = w(\pi)$, then the pair (x, π) satisfies the complementary slackness optimality conditions.*

Proof. Since $z(x) = w(\pi)$,

$$\sum_{(i,j) \in A} c_{ij}x_{ij} = \sum_{i \in N} b(i)\pi(i) - \sum_{(i,j) \in A} \max\{0, -c_{ij}^{\pi}\}u_{ij}. \quad (9.17)$$

Substituting the result of Property 2.4 in (9.17) shows that

$$\sum_{(i,j) \in A} \max\{0, -c_{ij}^{\pi}\}u_{ij} = \sum_{(i,j) \in A} -c_{ij}^{\pi}x_{ij}. \quad (9.18)$$

Now observe that both the sides have m terms, and each term on the left-hand side is nonnegative and its value is an upper bound on the corresponding term on the right-hand side (because $\max\{0, -c_{ij}^{\pi}\} \geq -c_{ij}^{\pi}$ and $u_{ij} \geq x_{ij}$). Therefore, the two sides can be equal only when

$$\max\{0, -c_{ij}^{\pi}\}u_{ij} = -c_{ij}^{\pi}x_{ij} \quad \text{for every arc } (i, j) \in A. \quad (9.19)$$

Now we consider three cases.

- (a) $c_{ij}^{\pi} > 0$. In this case, the left-hand side of (9.19) is zero, and the right-hand side can be zero only if $x_{ij} = 0$. This conclusion establishes (9.8a).
- (b) $0 < x_{ij} < u_{ij}$. In this case, $c_{ij}^{\pi} = 0$; otherwise, the right-hand side of (9.19) is negative. This conclusion establishes (9.8b).
- (c) $c_{ij}^{\pi} < 0$. In this case, the left-hand side of (9.19) is $-c_{ij}^{\pi}u_{ij}$ and therefore, $x_{ij} = u_{ij}$. This conclusion establishes (9.8c).

These results complete the proof of the theorem. \blacklozenge

The following result is an easy consequence of Theorems 9.6 and 9.7.

Property 9.8. *If x^* is an optimal solution of the minimum cost flow problem, and π is an optimal solution of the dual minimum cost flow problem, the pair (x^*, π) satisfies the complementary slackness optimality conditions (9.8).*

Proof. Theorem 9.6 implies that $z(x^*) = w(\pi)$ and Theorem 9.7 implies that the pair (x^*, π) satisfies (9.8). \blacklozenge

One important implication of the minimum cost flow duality is that it permits us to solve linear programs that have at most one $+1$ and at most one -1 in each row as minimum cost flow problems. Linear programs with this special structure arise in a variety of situations; Applications 19.10, 19.11, 19.18, and Exercises 9.9 and 19.18 provide a few examples.

Before examining the situation with at most one $+1$ and at most one -1 in each row, let us consider a linear program that has at most one $+1$ and at most one -1 in each column. We assume, without any loss of generality, that each constraint in the linear program is in equality form, because we can always bring the linear program into this form by introducing slack or surplus variables. (Observe that column corresponding the slack or surplus variables will also have one $+1$ or one -1 .) If each column has exactly one $+1$ and exactly one -1 , clearly the linear program is a minimum cost flow problem. Otherwise, we can augment this linear program by adding a redundant equality constraint which is the negative of the sum of all the

original constraints. (The new constraint corresponds to a new node that acts as a repository to deposit any excess supply or a source to fulfill any deficit demand from the other nodes.) The augmented linear program contains exactly one $+1$ and exactly one -1 in each column and the right-hand side values sum to zero. This model is clearly an instance of the minimum cost flow problem.

We now return to linear programs (in maximization form) that have at most one $+1$ and at most one -1 in each row. We allow a constraint in this linear program to be in any form: equality or inequality. The dual of this linear program contains at most one $+1$ and at most one -1 in each column, which we have already shown to be equivalent to a minimum cost flow problem. The variables in the dual problem will be nonnegative, nonpositive, or unrestricted, depending on whether they correspond to a less than or equal to, a greater than or equal to, or an equality constraint in the primal. A nonnegative variable x_{ij} defines a directed arc (i, j) in the resulting minimum cost flow formulation. To model any unrestricted variable x_{ij} , we replace it with two nonnegative variables, which is equivalent to introducing two arcs (i, j) and (j, i) of the same cost and capacity as this variable. The following theorem summarizes the preceding discussion.

Theorem 9.9. *Any linear program that contains (a) at most one $+1$ and at most one -1 in each column, or (b) at most one $+1$ and at most one -1 in each row, can be transformed into a minimum cost flow problem.* ♦

Minimum cost flow duality has several important implications. Since almost all algorithms for solving the primal problem also generate optimal node potentials $\pi(i)$ and the variables α_{ij} , solving the primal problem almost always solves both the primal and dual problems. Similarly, solving the dual problem typically solves the primal problem as well. Most algorithms for solving network flow problems explicitly or implicitly use properties of dual variables (since they are the node potentials that we have used at every turn) and of the dual linear program. In particular, the dual problem provides us with a certificate that if we can find a feasible dual solution that has the same objective function value as a given primal solution, we know from the strong duality theorem that the primal solution must be optimal, *without* making additional calculations and without considering other potentially optimal primal solutions. This certification procedure is a very powerful idea in network optimization, and in optimization in general. We have used it at many points in our previous developments and will see it many times again.

For network flow problems, the primal and dual problems are closely related via the basic shortest path and maximum flow problems that we have studied in previous chapters. In fact, these relationships help us to understand the fundamental importance of these two core problems to network flow theory and algorithms. We develop these relationships in the next section.

9.5 RELATING OPTIMAL FLOWS TO OPTIMAL NODE POTENTIALS

We next address the following questions: (1) Given an optimal flow, how might we obtain optimal node potentials? Conversely, (2) given optimal node potentials, how might we obtain an optimal flow? We show how to solve these problems by solving

either a shortest path problem or a maximum flow problem. These results point out an interesting relationship between the minimum cost flow problem and the maximum flow and shortest path problems.

Computing Optimal Node Potentials

We show that given an optimal flow x^* , we can obtain optimal node potentials by solving a shortest path problem (with possibly negative arc lengths). Let $G(x^*)$ denote the residual network with respect to the flow x^* . Clearly, $G(x^*)$ does not contain any negative cost cycle, for otherwise we would contradict the optimality of the solution x^* . Let $d(\cdot)$ denote the shortest path distances from node 1 to the rest of the nodes in the residual network if we use c_{ij} as arc lengths. The distances $d(\cdot)$ are well defined because the residual network does not contain a negative cycle. The shortest path optimality conditions (5.2) imply that

$$d(j) \leq d(i) + c_{ij} \quad \text{for all } (i, j) \text{ in } G(x^*). \quad (9.20)$$

Let $\pi = -d$. Then we can restate (9.20) as

$$c_{ij}^{\pi} = c_{ij} - \pi(i) + \pi(j) \geq 0 \quad \text{for all } (i, j) \text{ in } G(x^*).$$

Theorem 9.3 shows that π constitutes an optimal set of node potentials.

Obtaining Optimal Flows

We now show that given a set of optimal node potentials π , we can obtain an optimal solution x^* by solving a maximum flow problem. First, we compute the reduced cost c_{ij}^{π} of every arc $(i, j) \in A$ and then we examine all arcs one by one. We classify each arc (i, j) in one of the following ways and use these categorizations of the arcs to define a maximum flow problem.

Case 1: $c_{ij}^{\pi} > 0$

The condition (9.8a) implies that x_{ij}^* must be zero. We enforce this constraint by setting $x_{ij}^* = 0$ and deleting arc (i, j) from the network.

Case 2: $c_{ij}^{\pi} < 0$

The condition (9.8c) implies that $x_{ij}^* = u_{ij}$. We enforce this constraint by setting $x_{ij}^* = u_{ij}$ and deleting arc (i, j) from the network. Since we sent u_{ij} units of flow on arc (i, j) , we must decrease $b(i)$ by u_{ij} and increase $b(j)$ by u_{ij} .

Case 3: $c_{ij}^{\pi} = 0$

In this case we allow the flow on arc (i, j) to assume any value between 0 and u_{ij} .

Let $G' = (N, A')$ denote the resulting network and let b' denote the modified supplies/demands of the nodes. Now the problem reduces to finding a feasible flow in the network G' that meets the modified supplies/demands of the nodes. As noted in Section 6.2, we can find such a flow by solving a maximum flow problem defined as follows. We introduce a *source node* s , and a *sink node* t . For each node i with

$b'(i) > 0$, we add an arc (s, i) with capacity $b'(i)$ and for each node i with $b'(i) < 0$, we add an arc (i, t) with capacity $-b'(i)$. We now solve a maximum flow problem from node s to t in the transformed network obtaining a maximum flow x^* . The solution x_{ij}^* for all $(i, j) \in A$ is an optimal flow for the minimum cost flow problem in G .

9.6 CYCLE-CANCELING ALGORITHM AND THE INTEGRALITY PROPERTY

The negative cycle optimality conditions suggests one simple algorithmic approach for solving the minimum cost flow problem, which we call the *cycle-canceling algorithm*. This algorithm maintains a feasible solution and at every iteration attempts to improve its objective function value. The algorithm first establishes a feasible flow x in the network by solving a maximum flow problem (see Section 6.2). Then it iteratively finds negative cost-directed cycles in the residual network and augments flows on these cycles. The algorithm terminates when the residual network contains no negative cost-directed cycle. Theorem 9.1 implies that when the algorithm terminates, it has found a minimum cost flow. Figure 9.7 specifies this generic version of the cycle-canceling algorithm.

```

algorithm cycle-canceling;
begin
    establish a feasible flow  $x$  in the network;
    while  $G(x)$  contains a negative cycle do
        begin
            use some algorithm to identify a negative cycle  $W$ ;
             $\delta := \min\{r_{ij} : (i, j) \in W\}$ ;
            augment  $\delta$  units of flow in the cycle  $W$  and update  $G(x)$ ;
        end;
    end;

```

Figure 9.7 Cycle canceling algorithm.

We use the example shown in Figure 9.8(a) to illustrate the cycle-canceling algorithm. (The reader might notice that our example does not satisfy Assumption 9.4; we violate this assumption so that the network is simpler to analyze.) Figure 9.8(a) depicts a feasible flow in the network and Figure 9.8(b) gives the corresponding residual network. Suppose that the algorithm first selects the cycle 4–2–3–4 whose cost is -1 . The residual capacity of this cycle is 2. The algorithm augments 2 units of flow along this cycle. Figure 9.8(c) shows the modified residual network. In the next iteration, suppose that the algorithm selects the cycle 4–2–1–3–4 whose cost is -2 . The algorithm sends 1 unit of flow along this cycle. Figure 9.8(d) depicts the updated residual network. Since this residual network contains no negative cycle, the algorithm terminates.

In Chapter 5 we discussed several algorithms for identifying a negative cycle if one exists. One algorithm for identifying a negative cycle is the FIFO label-correcting algorithm for the shortest path problem described in Section 5.4; this algorithm requires $O(nm)$ time. We describe other algorithms for detecting negative cycles in Sections 11.7 and 12.7.

A by-product of the cycle-canceling algorithm is the following important result.

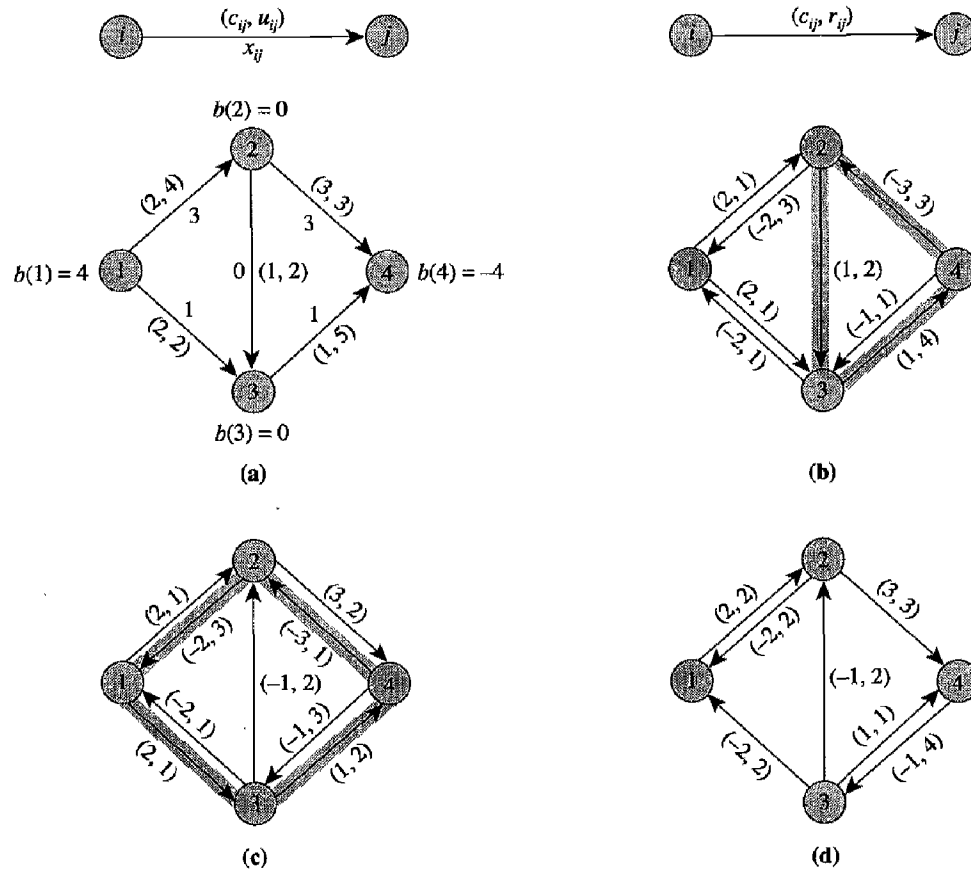


Figure 9.8 Illustrating the cycle canceling algorithm: (a) network example with a feasible flow x ; (b) residual network $G(x)$; (c) residual network after augmenting 2 units along the cycle 4-2-3-4; (d) residual network after augmenting 1 unit along the cycle 4-2-1-3-4.

Theorem 9.10 (Integrality Property). *If all arc capacities and supplies/demands of nodes are integer, the minimum cost flow problem always has an integer minimum cost flow.*

Proof. We show this result by performing induction on the number of iterations. The algorithm first establishes a feasible flow in the network by solving a maximum flow problem. By Theorem 6.5 the problem has an integer feasible flow and we assume that the maximum flow algorithm finds an integer solution since all arc capacities in the network are integer and the initial residual capacities are also integer. The flow augmented by the cycle-canceling algorithm in any iteration equals the minimum residual capacity in the cycle canceled, which by the inductive hypothesis is integer. Therefore the modified residual capacities in the next iteration will again be integer. This conclusion implies the assertion of the theorem. ♦

Let us now consider the number of iterations that the algorithm performs. For the minimum cost flow problem, mCU is an upper bound on the initial flow cost

[since $c_{ij} \leq C$ and $x_{ij} \leq U$ for all $(i, j) \in A$] and $-mCU$ is a lower bound on the optimal flow cost [since $c_{ij} \geq -C$ and $x_{ij} \leq U$ for all $(i, j) \in A$]. Each iteration of the cycle-canceling algorithm changes the objective function value by an amount $(\sum_{(i,j) \in W} c_{ij})\delta$, which is strictly negative. Since we are assuming that all the data of the problem are integral, the algorithm terminates within $O(mCU)$ iterations and runs in $O(nm^2CU)$ time.

The generic version of the cycle-canceling algorithm does not specify the order for selecting negative cycles from the network. Different rules for selecting negative cycles produce different versions of the algorithm, each with different worst-case and theoretical behavior. The network simplex algorithm, which is widely considered to be one of the fastest algorithms for solving the minimum cost flow problem in practice, is a particular version of the cycle-canceling algorithm. The network simplex algorithm maintains information (a spanning tree solution and node potentials) that enables it to identify a negative cost cycle in $O(m)$ time. However, due to *degeneracy*, the algorithm cannot necessarily send a positive amount of flow along this cycle. We discuss these issues in Chapter 11, where we consider the network simplex algorithm in more detail. The most general implementation of the network simplex algorithm does not run in polynomial time. The following two versions of the cycle-canceling algorithm are, however, polynomial-time implementations.

Augmenting flow in a negative cycle with maximum improvement.

Let x be any feasible flow and let x^* be an optimal flow. The improvement in the objective function value due to an augmentation along a cycle W is $-(\sum_{(i,j) \in W} c_{ij}) (\min\{r_{ij} : (i, j) \in W\})$. We observed in the proof of Theorem 3.7 in Section 3.5 that x^* equals x plus the flow on at most m augmenting cycles with respect to x , and improvements in cost due to flow augmentations on these augmenting cycles sum to $cx - cx^*$. Consequently, at least one of these augmenting cycles with respect to x must decrease the objective function value by at least $(cx - cx^*)/m$. Consequently, if the algorithm always augments flow along a cycle giving the maximum possible improvement, then Theorem 3.1 implies that the method would obtain an optimal flow within $O(m \log(mCU))$ iterations. Finding a maximum improvement cycle is difficult (i.e., it is a \mathcal{NP} -complete problem), but a modest variation of this approach yields a polynomial-time algorithm for the minimum cost flow problem. We provide a reference for this algorithm in the reference notes.

Augmenting flow along a negative cycle with minimum mean cost.

We define the *mean cost* of a cycle as its cost divided by the number of arcs it contains. A *minimum mean cycle* is a cycle whose mean cost is as small as possible. It is possible to identify a minimum mean cycle in $O(nm)$ or $O(\sqrt{n} m \log(nC))$ time (see the reference notes of Chapter 5). Researchers have shown that if the cycle-canceling algorithm always augments flow along a minimum mean cycle, it performs $O(\min\{nm \log(nC), nm^2 \log n\})$ iterations. We describe this algorithm in Section 10.5.

9.7 SUCCESSIVE SHORTEST PATH ALGORITHM

The cycle-canceling algorithm maintains feasibility of the solution at every step and attempts to achieve optimality. In contrast, the successive shortest path algorithm maintains optimality of the solution (as defined in Theorem 9.3) at every step and strives to attain feasibility. It maintains a solution x that satisfies the nonnegativity and capacity constraints, but violates the mass balance constraints of the nodes. At each step, the algorithm selects a node s with excess supply (i.e., supply not yet sent to some demand node) and a node t with unfulfilled demand and sends flow from s to t along a shortest path in the residual network. The algorithm terminates when the current solution satisfies all the mass balance constraints.

To describe this algorithm as well as several later developments, we first introduce the concept of *pseudoflows*. A *pseudoflow* is a function $x: A \rightarrow \mathbb{R}^+$ satisfying only the capacity and nonnegativity constraints; it need not satisfy the mass balance constraints. For any pseudoflow x , we define the *imbalance* of node i as

$$e(i) = b(i) + \sum_{\{j: (j,i) \in A\}} x_{ji} - \sum_{\{j: (i,j) \in A\}} x_{ij} \quad \text{for all } i \in N.$$

If $e(i) > 0$ for some node i , we refer to $e(i)$ as the *excess* of node i ; if $e(i) < 0$, we call $-e(i)$ the node's *deficit*. We refer to a node i with $e(i) = 0$ as *balanced*. Let E and D denote the sets of excess and deficit nodes in the network. Notice that $\sum_{i \in N} e(i) = \sum_{i \in N} b(i) = 0$, and hence $\sum_{i \in E} e(i) = -\sum_{i \in D} e(i)$. Consequently, if the network contains an excess node, it must also contain a deficit node. The residual network corresponding to a pseudoflow is defined in the same way that we define the residual network for a flow.

Using the concept of pseudoflow and the reduced cost optimality conditions specified in Theorem 9.3, we next prove some results that we will use extensively in this and the following chapters.

Lemma 9.11. *Suppose that a pseudoflow (or a flow) x satisfies the reduced cost optimality conditions with respect to some node potentials π . Let the vector d represent the shortest path distances from some node s to all other nodes in the residual network $G(x)$ with c_{ij}^{π} as the length of an arc (i, j) . Then the following properties are valid:*

- (a) *The pseudoflow x also satisfies the reduced cost optimality conditions with respect to the node potentials $\pi' = \pi - d$.*
- (b) *The reduced costs $c_{ij}^{\pi'}$ are zero for all arcs (i, j) in a shortest path from node s to every other node.*

Proof. Since x satisfies the reduced cost optimality conditions with respect to π , $c_{ij}^{\pi} \geq 0$ for every arc (i, j) in $G(x)$. Furthermore, since the vector d represents shortest path distances with c_{ij}^{π} as arc lengths, it satisfies the shortest path optimality conditions, that is,

$$d(j) \leq d(i) + c_{ij}^{\pi} \quad \text{for all } (i, j) \text{ in } G(x). \quad (9.21)$$

Substituting $c_{ij}^{\pi'} = c_{ij} - \pi(i) + \pi(j)$ in (9.21), we obtain $d(j) \leq d(i) + c_{ij} - \pi(i) + \pi(j)$. Alternatively, $c_{ij} - (\pi(i) - d(i)) + (\pi(j) - d(j)) \geq 0$, or $c_{ij}^{\pi'} \geq 0$. This conclusion establishes part (a) of the lemma.

Consider next a shortest path from node s to some node l . For each arc (i, j) in this path, $d(j) = d(i) + c_{ij}^{\pi}$. Substituting $c_{ij}^{\pi} = c_{ij} - \pi(i) + \pi(j)$ in this equation, we obtain $c_{ij}^{\pi} = 0$. This conclusion establishes part (b) of the lemma. ♦

The following result is an immediate corollary of the preceding lemma.

Lemma 9.12. *Suppose that a pseudoflow (or a flow) x satisfies the reduced cost optimality conditions and we obtain x' from x by sending flow along a shortest path from node s to some other node k ; then x' also satisfies the reduced cost optimality conditions.*

Proof. Define the potentials π and π' as in Lemma 9.11. The proof of Lemma 9.11 implies that for every arc (i, j) in the shortest path P from node s to the node k , $c_{ij}^{\pi} = 0$. Augmenting flow on any such arc might add its reversal (j, i) to the residual network. But since $c_{ij}^{\pi} = 0$ for each arc $(i, j) \in P$, $c_{ji}^{\pi} = 0$ and the arc (j, i) also satisfies the reduced cost optimality conditions. These results establish the lemma. ♦

We are now in a position to describe the successive shortest path algorithm. The node potentials play a very important role in this algorithm. Besides using them to prove the correctness of the algorithm, we use them to maintain nonnegative arc lengths so that we can solve the shortest path problem more efficiently. Figure 9.9 gives a formal statement of the successive shortest path algorithm.

We illustrate the successive shortest path algorithm on the same numerical example we used to illustrate the cycle canceling algorithm. Figure 9.10(a) shows the initial residual network. Initially, $E = \{1\}$ and $D = \{4\}$. Therefore, in the first iteration, $s = 1$ and $t = 4$. The shortest path distances d (with respect to the reduced costs) are $d = (0, 2, 2, 3)$ and the shortest path from node 1 to node 4 is 1–3–4. Figure 9.10(b) shows the updated node potentials and reduced costs, and Figure 9.10(c) shows the solution after we have augmented $\min\{e(1), -e(4), r_{13}, r_{34}\} = \min\{4, 4, 2, 5\} = 2$ units of flow along the path 1–3–4. In the second iteration, $k =$

```

algorithm successive shortest path;
begin
   $x := 0$  and  $\pi := 0$ ;
   $e(i) := b(i)$  for all  $i \in N$ ;
  initialize the sets  $E := \{i : e(i) > 0\}$  and  $D := \{i : e(i) < 0\}$ ;
  while  $E \neq \emptyset$  do
    begin
      select a node  $k \in E$  and a node  $l \in D$ ;
      determine shortest path distances  $d(j)$  from node  $s$  to all
        other nodes in  $G(x)$  with respect to the reduced costs  $c_{ij}^{\pi}$ ;
      let  $P$  denote a shortest path from node  $k$  to node  $l$ ;
      update  $\pi := \pi - d$ ;
       $\delta := \min\{e(k), -e(l), \min\{r_{ij} : (i, j) \in P\}\}$ ;
      augment  $\delta$  units of flow along the path  $P$ ;
      update  $x$ ,  $G(x)$ ,  $E$ ,  $D$ , and the reduced costs;
    end;
  end;

```

Figure 9.9 Successive shortest path algorithm.

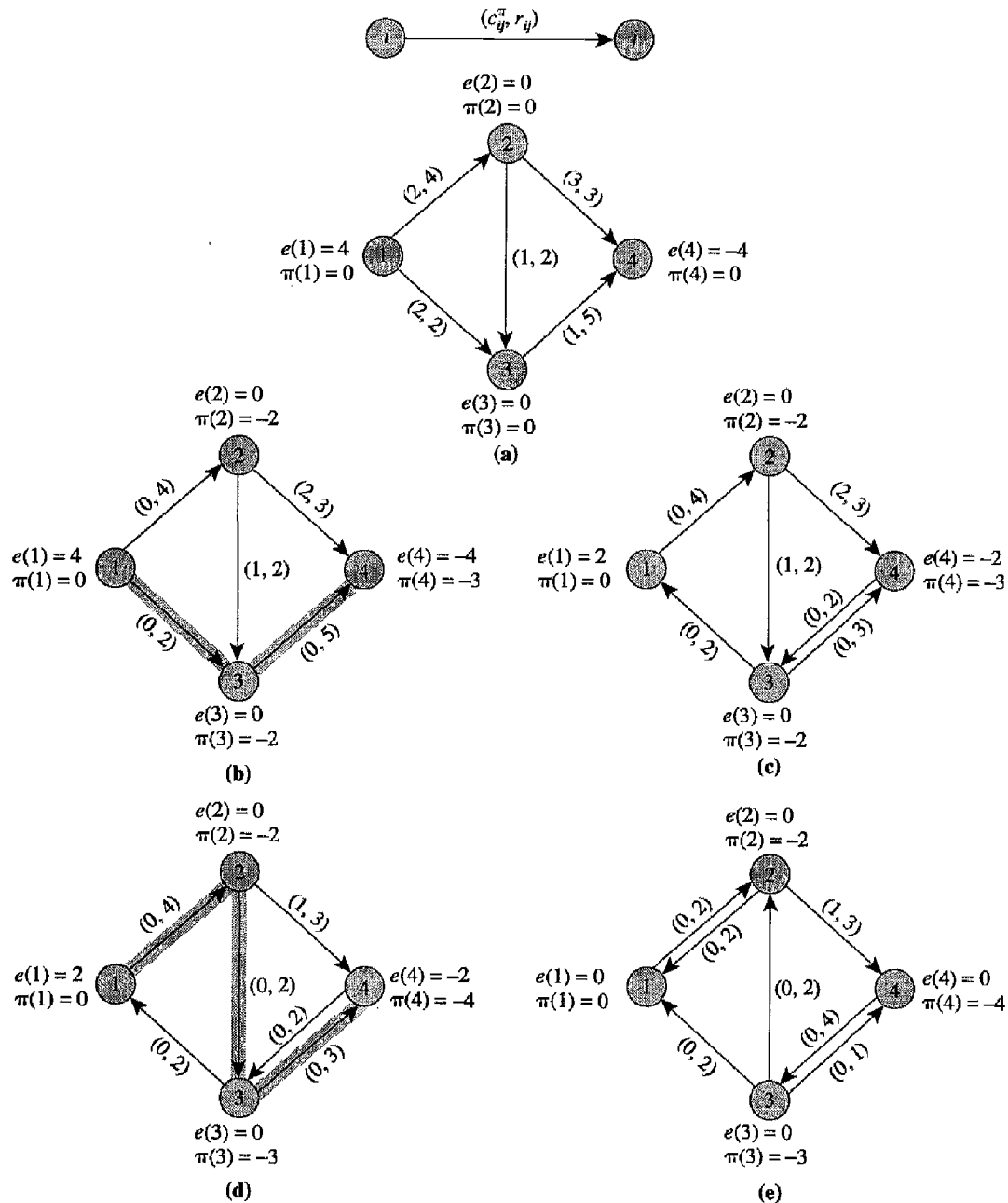


Figure 9.10 Illustrating the successive shortest path algorithm: (a) initial residual network for $x = 0$ and $\pi = 0$; (b) network after updating the potentials π ; (c) network after augmenting 2 units along the path 1-3-4; (d) network after updating the potentials π ; (e) network after augmenting 2 units along the path 1-2-3-4.

1, $l = 4$, $d = (0, 0, 1, 1)$ and the shortest path from node 1 to node 4 is 1-2-3-4. Figure 9.10(d) shows the updated node potentials and reduced costs, and Figure 9.10(e) shows the solution after we have augmented $\min\{e(1), -e(4), r_{12}, r_{23}, r_{34}\} = \min\{2, 2, 4, 2, 3\} = 2$ units of flow. At the end of this iteration, all imbalances become zero and the algorithm terminates.

We now justify the successive shortest path algorithm. To initialize the algorithm, we set $x = 0$, which is a feasible pseudoflow. For the zero pseudoflow x , $G(x) = G$. Note that this solution together with $\pi = 0$ satisfies the reduced cost optimality conditions because $c_{ij}^{\pi} = c_{ij} \geq 0$ for every arc (i, j) in the residual network $G(x)$ (recall Assumption 9.5, which states that all arc costs are nonnegative). Observe that as long as any node has a nonzero imbalance, both E and D must be nonempty since the total sum of excesses equals the total sum of deficits. Thus until all nodes are balanced, the algorithm always succeeds in identifying an excess node k and a deficit node l . Assumption 9.4 implies that the residual network contains a directed path from node k to every other node, including node l . Therefore, the shortest path distances $d(\cdot)$ are well defined. Each iteration of the algorithm solves a shortest path problem with nonnegative arc lengths and strictly decreases the excess of some node (and, also, the deficit of some other node). Consequently, if U is an upper bound on the largest supply of any node, the algorithm would terminate in at most nU iterations. If $S(n, m, C)$ denotes the time taken to solve a shortest path problem with nonnegative arc lengths, the overall complexity of this algorithm is $O(nUS(n, m, nC))$. [Note that we have used nC rather than C in this expression, since the costs in the residual network are bounded by nC .] We refer the reader to the reference notes of Chapter 4 for the best available value of $S(n, m, C)$.

The successive shortest path algorithm requires pseudopolynomial time to solve the minimum cost flow problem since it is polynomial in n , m and the largest supply U . This algorithm is, however, polynomial time for the assignment problem, a special case of the minimum cost flow problem, for which $U = 1$. In Chapter 10, using scaling techniques, we develop weakly and strongly polynomial-time versions of the successive shortest path algorithm. In Section 14.5 we generalize this approach even further, developing a polynomial-time algorithm for the convex cost flow problem.

We now suggest some practical improvements to the successive shortest path algorithm. As stated, this algorithm selects an excess node k , uses Dijkstra's algorithm to identify shortest paths from node k to all other nodes, and augments flow along a shortest path from node k to some deficit node l . In fact, it is not necessary to determine a shortest path from node k to *all* nodes; a shortest path from node k to *one* deficit node l is sufficient. Consequently, we could terminate Dijkstra's algorithm whenever it permanently labels the first deficit node l . At this point we might modify the node potentials in the following manner:

$$\pi(i) = \begin{cases} \pi(i) - d(i) & \text{if node } i \text{ is permanently labeled} \\ \pi(i) - d(l) & \text{if node } i \text{ is temporarily labeled.} \end{cases}$$

In Exercise 9.47 we ask the reader to show that with this choice of the modified node potentials, the reduced costs of all the arcs in the residual network remain nonnegative and the reduced costs of the arcs along the shortest path from node k to node l are zero. Observe that we can alternatively modify the node potentials in the following manner:

$$\pi(i) = \begin{cases} \pi(i) - d(i) + d(l) & \text{if node } i \text{ is permanently labeled} \\ \pi(i) & \text{if node } i \text{ is temporarily labeled.} \end{cases}$$

This scheme for updating node potentials is the same as the previous scheme except that we add $d(l)$ to all of the node potentials (which does not affect the reduced cost of any arc). An advantage of this scheme is that the algorithm spends no time updating the potentials of the temporarily labeled nodes.

9.8 PRIMAL-DUAL ALGORITHM

The primal-dual algorithm for the minimum cost flow problem is similar to the successive shortest path algorithm in the sense that it also maintains a pseudoflow that satisfies the reduced cost optimality conditions and gradually converts it into a flow by augmenting flows along shortest paths. In contrast, instead of sending flow along one shortest path at a time, it solves a maximum flow problem that sends flow along all shortest paths.

The primal-dual algorithm generally transforms the minimum cost flow problem into a problem with a single excess node and a single deficit node. We transform the problem into this form by introducing a *source* node s and a *sink* node t . For each node i with $b(i) > 0$, we add a zero cost arc (s, i) with capacity $b(i)$, and for each node i with $b(i) < 0$, we add a zero cost arc (i, t) with capacity $-b(i)$. Finally, we set $b(s) = \sum_{i \in N: b(i) > 0} b(i)$, $b(t) = -b(s)$, and $b(i) = 0$ for all $i \in N$. It is easy to see that a minimum cost flow in the transformed network gives a minimum cost flow in the original network. For simplicity of notation, we shall represent the transformed network as $G = (N, A)$, which is the same representation that we used for the original network.

The primal-dual algorithm solves a maximum flow problem on a subgraph of the residual network $G(x)$, called the *admissible network*, which we represent as $G^o(x)$. We define the admissible network $G^o(x)$ with respect to a pseudoflow x that satisfies the reduced cost optimality conditions for some node potentials π ; the admissible network contains only those arcs in $G(x)$ with a zero reduced cost. The residual capacity of an arc in $G^o(x)$ is the same as that in $G(x)$. Observe that every directed path from node s to node t in $G^o(x)$ is a shortest path in $G(x)$ between the same pair of nodes (see Exercise 5.20). Figure 9.11 formally describes the primal-dual algorithm on the transformed network.

```

algorithm primal-dual;
begin
   $x := 0$  and  $\pi := 0$ ;
   $e(s) := b(s)$  and  $e(t) := b(t)$ ;
  while  $e(s) > 0$  do
    begin
      determine shortest path distances  $d(\cdot)$  from node  $s$  to all other nodes in  $G(x)$  with
        respect to the reduced costs  $c_j$ ;
      update  $\pi := \pi - d$ ;
      define the admissible network  $G^o(x)$ ;
      establish a maximum flow from node  $s$  to node  $t$  in  $G^o(x)$ ;
      update  $e(s)$ ,  $e(t)$ , and  $G(x)$ ;
    end;
  end;

```

Figure 9.11 Primal-dual algorithm.

To illustrate the primal–dual algorithm, we consider the numerical example shown in Figure 9.12(a). Figure 9.12(b) shows the transformed network. The shortest path computation yields the vector $d = (0, 0, 0, 1, 2, 1)$ whose components are in

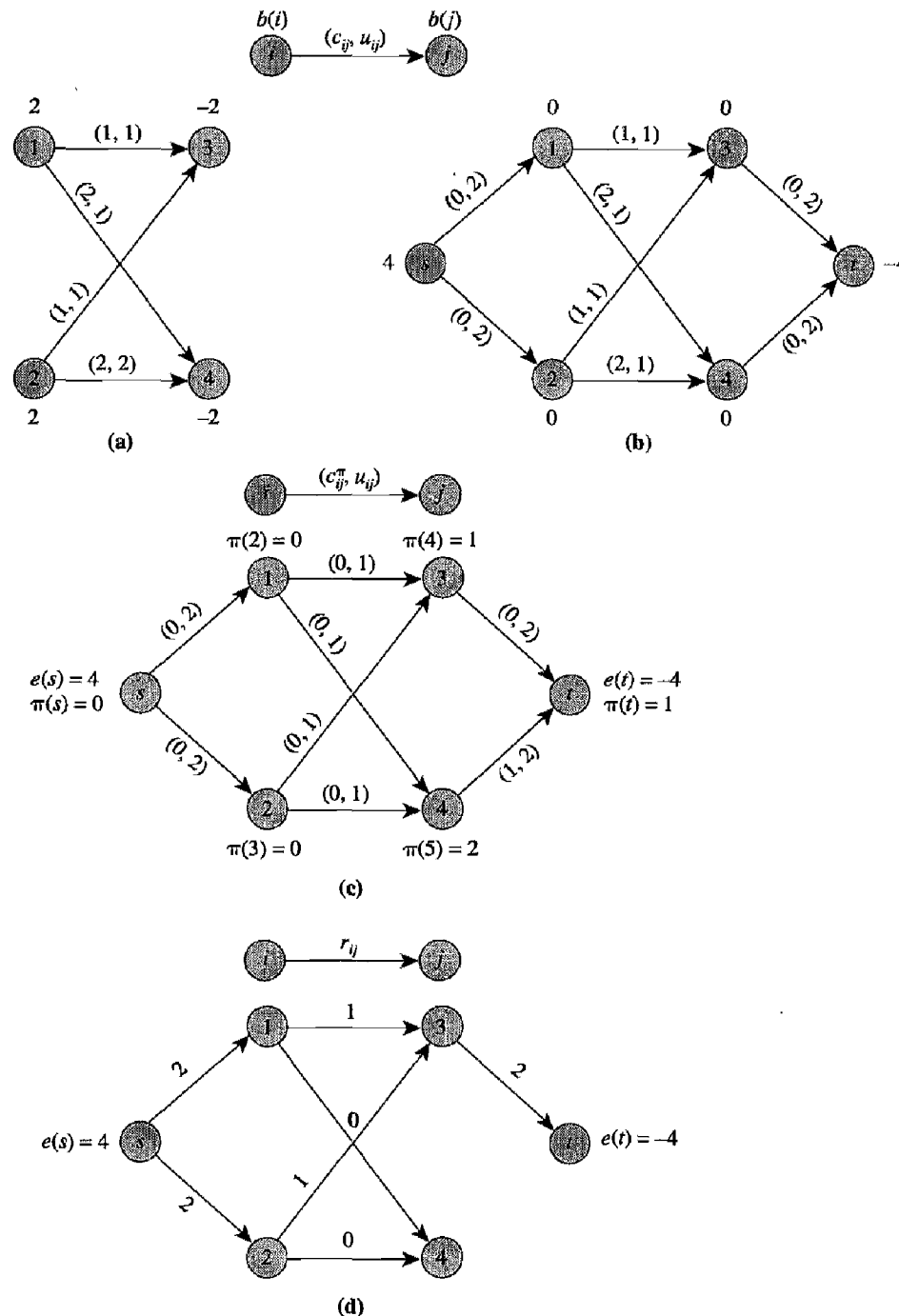


Figure 9.12 Illustrating the primal–dual algorithm: (a) example network; (b) transformed network; (c) residual network after updating the node potentials; (d) admissible network.

the order $s, 1, 2, 3, 4, t$. Figure 9.12(c) shows the modified node potentials and reduced costs and Figure 9.12(d) shows the admissible network at this stage in the computations. When we apply the maximum flow algorithm to the admissible network, it is able to send 2 units of flow from node s to node t . Observe that the admissible network contained two paths from node s to node t and the maximum flow computation saturates both the paths. The successive shortest path algorithm would have taken two iterations to send the 2 units of flow. As the reader can verify, the second iteration of the primal–dual algorithm also sends 2 units of flow from node s to node t , at which point it converts the pseudoflow into a flow and terminates.

The primal–dual algorithm guarantees that the excess of node s strictly decreases at each iteration, and also assures that the node potential of the sink strictly decreases from one iteration to the next. The second observation follows from the fact that once we have established a maximum flow in $G^o(x)$, the residual network $G(x)$ contains no directed path from node s to node t consisting entirely of arcs with zero reduced costs. Consequently, in the next iteration, when we solve the shortest path problem, $d(t) \geq 1$. These observations give a bound of $\min\{nU, nC\}$ on the number of iterations since initially $e(s) \leq nU$, and the value of no node potential can fall below $-nC$ (see Exercise 9.25). This bound on the number of iterations is better than that of the successive shortest path algorithm, but, of course, the algorithm incurs the additional expense of solving a maximum flow problem at every iteration. If $S(n, m, C)$ and $M(n, m, U)$ denote the solution times of shortest path and the maximum flow algorithms, the primal–dual algorithm has an overall complexity of $O(\min\{nU, nC\} \cdot \{S(n, m, nC) + M(n, m, U)\})$.

In concluding this discussion, we might comment on why this algorithm is known as the primal–dual algorithm. This name stems from linear programming duality theory. In the linear programming literature, the primal–dual algorithm always maintains a dual feasible solution π and a primal solution that might violate some supply/demand constraints (i.e., is primal infeasible), so that the pair satisfies the complementary slackness conditions. For a given dual feasible solution, the algorithm attempts to decrease the degree of primal infeasibility to the minimum possible level. [Recall that the algorithm solves a maximum flow problem to reduce $e(s)$ by the maximum amount.] When no further reduction in the primal infeasibility is possible, the algorithm modifies the dual solution (i.e., node potentials in the network flow context) and again tries to minimize primal infeasibility. (This primal–dual approach is applicable to several combinatorial optimization problems and also to the general linear programming problem. Indeed, this primal–dual solution strategy is one of the most popular approaches for solving specially structured problems and has often yielded fairly efficient and intuitively appealing algorithms.)

9.9 OUT-OF-KILTER ALGORITHM

The successive shortest path and primal–dual algorithms maintain a solution that satisfies the reduced cost optimality conditions and the flow bound constraints but violates the mass balance constraints. These algorithms iteratively modify arc flows and node potentials so that the flow at each step comes closer to satisfying the mass balance constraints. However, we could just as well have developed other solution strategies by violating other constraints at intermediate steps. The out-of-kilter al-

gorithm, which we discuss in this section, satisfies only the mass balance constraints, so intermediate solutions might violate both the optimality conditions and the flow bound restrictions. The algorithm iteratively modifies flows and potentials in a way that decreases the infeasibility of the solution (in a way to be specified) and, simultaneously, moves it closer to optimality. In essence, the out-of-kilter algorithm is similar to the successive shortest path and primal–dual algorithms because its fundamental step at every iteration is solving a shortest path problem and augmenting flow along a shortest path.

To describe the out-of-kilter algorithm, we refer to the complementary slackness optimality conditions stated in Theorem 9.4. For ease of reference, let us restate these conditions.

$$\text{If } x_{ij} = 0, \text{ then } c_{ij}^{\pi} \geq 0. \quad (9.22a)$$

$$\text{If } 0 < x_{ij} < u_{ij}, \text{ then } c_{ij}^{\pi} = 0. \quad (9.22b)$$

$$\text{If } x_{ij} = u_{ij}, \text{ then } c_{ij}^{\pi} \leq 0, \quad (9.22c)$$

The name *out-of-kilter algorithm* reflects the fact that arcs in the network either satisfy the complementary slackness optimality conditions (are *in-kilter*) or do not (are *out-of-kilter*). The so-called *kilter diagram* is a convenient way to represent these conditions. As shown in Figure 9.13, the kilter diagram of an arc (i, j) is the collection of all points (x_{ij}, c_{ij}^{π}) in the two-dimensional plane that satisfy the optimality conditions (9.22). The condition 9.22(a) implies that $c_{ij}^{\pi} \geq 0$ if $x_{ij} = 0$; therefore, the kilter diagram contains all points with zero x_{ij} -coordinates and nonnegative c_{ij}^{π} -coordinates. Similarly, the condition 9.22(b) yields the horizontal segment of the diagram, and condition 9.22(c) yields the other vertical segment of the diagram. Each arc has its own kilter diagram.

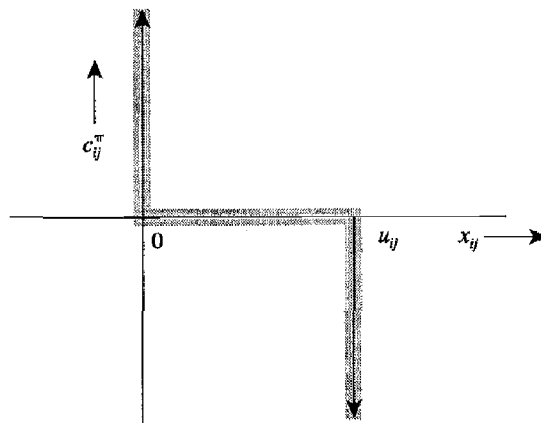


Figure 9.13 Kilter diagram for arc (i, j) .

Notice that for every arc (i, j) , the flow x_{ij} and reduced cost c_{ij}^{π} define a point (x_{ij}, c_{ij}^{π}) in the two-dimensional plane. If the point (x_{ij}, c_{ij}^{π}) lies on the thick lines in the kilter diagram, the arc is *in-kilter*; otherwise, it is *out-of-kilter*. For instance, the points B , D , and E in Figure 9.14 are *in-kilter*, whereas the points A and C are *out-of-kilter*. We define the *kilter number* k_{ij} of each arc (i, j) in A as the magnitude of the change in x_{ij} required to make the arc an *in-kilter* arc while keeping c_{ij}^{π} fixed.

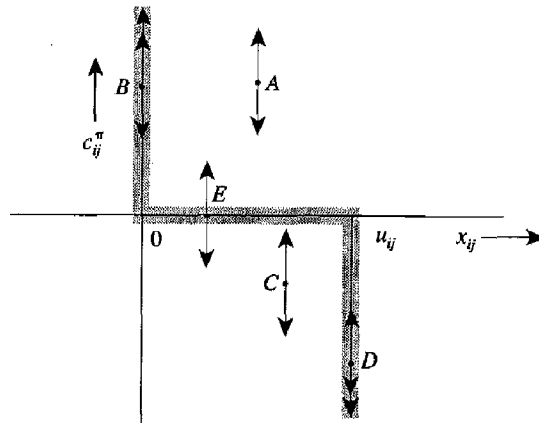


Figure 9.14 Examples of in-kilter and out-of-kilter arcs.

Therefore, in accordance with conditions (9.22a) and (9.22c), if $c_{ij}^\pi > 0$, then $k_{ij} = |x_{ij}|$, and if $c_{ij}^\pi < 0$, then $k_{ij} = |u_{ij} - x_{ij}|$. If $c_{ij}^\pi = 0$ and $x_{ij} > u_{ij}$, then $k_{ij} = x_{ij} - u_{ij}$. If $c_{ij}^\pi = 0$ and $x_{ij} < 0$, then $k_{ij} = -x_{ij}$. The kilter number of any in-kilter arc is zero. The sum $K = \sum_{(i,j) \in A} k_{ij}$ of all kilter numbers provides us with a measure of how far the current solution is from optimality; the smaller the value of K , the closer the current solution is to being an optimal solution.

In describing the out-of-kilter algorithm, we begin by making a simplifying assumption that the algorithm starts with a feasible flow. At the end of this section we show how to extend the algorithm so that it applies to situations when the initial flow does not satisfy the arc flow bounds (we also consider situations with nonzero lower bounds on arc flows).

To describe the out-of-kilter algorithm, we will work on the residual network; in this setting, the algorithm iteratively decreases the kilter number of one or more arcs in the residual network. To do so, we must be able to define the kilter number of the arcs in the residual network $G(x)$. We set the kilter number k_{ij} of an arc (i, j) in the following manner:

$$k_{ij} = \begin{cases} 0 & \text{if } c_{ij}^\pi \geq 0. \\ r_{ij} & \text{if } c_{ij}^\pi < 0. \end{cases} \quad (9.23)$$

This definition of the kilter number of an arc in the residual network is consistent with our previous definition: It is the change in flow (or, equivalently, the residual capacity) required so that the arc satisfies its optimality condition [which, in the case of residual networks, is the reduced cost optimality condition (9.7)]. An arc (i, j) in the residual network with $c_{ij}^\pi \geq 0$ satisfies its optimality condition (9.7), but an arc (i, j) with $c_{ij}^\pi < 0$ does not. In the latter case, we must send r_{ij} units of flow on the arc (i, j) so that it drops out of the residual network and thus satisfies its optimality condition.

The out-of-kilter algorithm maintains a feasible flow x and a set of node potentials π . We could obtain a feasible flow by solving a maximum flow problem (as described in Section 6.2) and start with $\pi = 0$. Subsequently, the algorithm maintains all of the in-kilter arcs as in-kilter arcs and successively transforms the out-of-kilter arcs into in-kilter arcs. The algorithm terminates when all arcs in the residual network become in-kilter. Figure 9.15 gives a formal description of the out-of-kilter algorithm.

```

algorithm out-of-kilter;
begin
   $\pi := 0$ ;
  establish a feasible flow  $x$  in the network;
  define the residual network  $G(x)$  and compute the kilter numbers of arcs;
  while the network contains an out-of-kilter arc do
    begin
      select an out-of-kilter arc  $(p, q)$  in  $G(x)$ ;
      define the length of each arc  $(i, j)$  in  $G(x)$  as  $\max\{0, c_{ij}^\pi\}$ ;
      let  $d(\cdot)$  denote the shortest path distances from node  $q$  to all other nodes in
         $G(x) - \{(q, p)\}$  and let  $P$  denote a shortest path from node  $q$  to node  $p$ ;
      update  $\pi'(i) := \pi(i) - d(i)$  for all  $i \in N$ ;
      if  $c_{pq}^{\pi'} < 0$  then
        begin
           $W := P \cup \{(p, q)\}$ ;
           $\delta := \min\{r_{ij} : (i, j) \in W\}$ ;
          augment  $\delta$  units of flow along  $W$ ;
          update  $x$ ,  $G(x)$ , and the reduced costs;
        end;
      end;
    end;
  end;

```

Figure 9.15 Out-of-kilter algorithm.

We now discuss the correctness and complexity of the out-of-kilter algorithm. The correctness argument of the algorithm uses the fact that kilter numbers of arcs are nonincreasing. Two operations in the algorithm affect the kilter numbers of arcs: updating node potentials and augmenting flow along the cycle W . In the next two lemmas we show that these operations do not increase the kilter number of any arc.

Lemma 9.13. *Updating the node potentials does not increase the kilter number of any arc in the residual network.*

Proof. Let π and π' denote the node potentials in the out-of-kilter algorithm before and after the update. The definition of the kilter numbers from (9.23) implies that the kilter number of an arc (i, j) can increase only if $c_{ij}^\pi \geq 0$ and $c_{ij}^{\pi'} < 0$. We show that this cannot happen. Consider any arc (i, j) with $c_{ij}^\pi \geq 0$. We wish to show that $c_{ij}^{\pi'} \geq 0$. Since $c_{pq}^{\pi'} < 0$, $(i, j) \neq (p, q)$. Since the distances $d(\cdot)$ represent the shortest path distances with $\max\{0, c_{ij}^\pi\}$ as the length of arc (i, j) , the shortest path distances satisfy the following shortest path optimality condition (see Section 5.2):

$$d(j) \leq d(i) + \max\{0, c_{ij}^\pi\} = d(i) + c_{ij}^\pi.$$

The equality in this expression is valid because, by assumption, $c_{ij}^\pi \geq 0$. The preceding expression shows that

$$c_{ij}^\pi + d(i) - d(j) = c_{ij}^{\pi'} \geq 0,$$

so each arc in the residual network with a nonnegative reduced cost has a nonnegative reduced cost after the potentials update, which implies the conclusion of the lemma. ♦

Lemma 9.14. *Augmenting flow along the directed cycle $W = P \cup \{(p, q)\}$ does not increase the kilter number of any arc in the residual network and strictly decreases the kilter number of the arc (p, q) .*

Proof. Notice that the flow augmentation can change the kilter number of only the arcs in $W = P \cup \{(p, q)\}$ and their reversals. Since P is a shortest path in the residual network with $\max\{0, c_{ij}^\pi\}$ as the length of arc (i, j) ,

$$d(j) = d(i) + \max\{0, c_{ij}^\pi\} \geq d(i) + c_{ij}^\pi \quad \text{for each arc } (i, j) \in P,$$

which, using $\pi' = \pi - d$ and the definition $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$, implies that

$$c_{ij}^{\pi'} \leq 0 \quad \text{for each arc } (i, j) \in P.$$

Since the reduced cost of each arc (i, j) in P with respect to π' is nonpositive, the condition (9.23) shows that sending additional flow does not increase the arc's kilter number, but might decrease it. The flow augmentation might add the reversals of arcs in P , but since $c_{ij}^{\pi'} \leq 0$, the reversal of this arc (j, i) has $c_{ji}^{\pi'} \geq 0$, and therefore arc (j, i) is an in-kilter arc.

Finally, we consider arc (p, q) . Recall from the algorithm description in Figure 9.15 that we augment flow along the arc (p, q) only if it is an out-of-kilter arc (i.e., $c_{pq}^{\pi'} < 0$). Since augmenting flow along the arc (p, q) decreases its residual capacity, the augmentation decreases this arc's kilter number. Since $c_{qp}^{\pi'} > 0$, arc (q, p) remains an in-kilter arc. These conclusions complete the proof of the lemma. ♦

The preceding two lemmas allow us to obtain a pseudopolynomial bound on the running time of the out-of-kilter algorithm. Initially, the kilter number of an arc is at most U ; therefore, the sum of the kilter numbers is at most mU . At each iteration, the algorithm selects an arc, say (p, q) , with a positive kilter number and either makes it an in-kilter arc during the potential update step or decreases its kilter number by the subsequent flow augmentation. Therefore, the sum of kilter numbers decreases by at least 1 unit at every iteration. Consequently, the algorithm terminates within $O(mU)$ iterations. The dominant computation within each iteration is solving a shortest path problem. Therefore, if $S(n, m, C)$ is the time required to solve a shortest path problem with nonnegative arc lengths, the out-of-kilter algorithm runs in $O(mU S(n, m, nC))$ time.

How might we modify the algorithm to handle situations when the arc flows do not necessarily satisfy their flow bounds? In examining this case we consider the more general problem setting by allowing the arcs to have nonzero lower bounds. Let l_{ij} denote the lower bound on the flow on arc $(i, j) \in A$. In this case, the complementary slackness optimality conditions become:

$$\text{If } x_{ij} = l_{ij}, \text{ then } c_{ij}^\pi \geq 0. \quad (9.24a)$$

$$\text{If } l_{ij} < x_{ij} < u_{ij}, \text{ then } c_{ij}^\pi = 0. \quad (9.24b)$$

$$\text{If } x_{ij} = u_{ij}, \text{ then } c_{ij}^\pi \leq 0. \quad (9.24c)$$

The thick lines in Figure 9.16 define the kilter diagram for this case. Consider arc (i, j) . If the point (x_{ij}, c_{ij}^π) lies on the thick line in Figure 9.16, the arc is an in-kilter arc; otherwise it is an out-of-kilter arc. As earlier, we define the kilter number

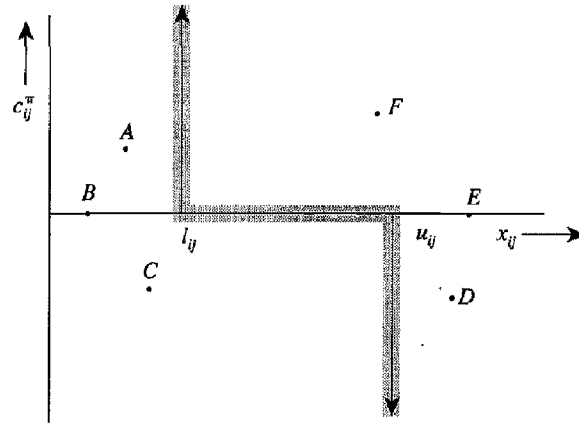


Figure 9.16 Kilter diagram for an arc (i, j) with a nonzero lower bound.

of an arc (i, j) in A as the magnitude of the change in x_{ij} required to make the arc an in-kilter arc while keeping c_{ij}^π fixed. Since arcs might violate their flow bounds, six types of out-of-kilter arcs are possible, which we depict by points A , B , C , D , E , and F in Figure 9.16. For example, the kilter numbers of arcs with coordinates depicted by the points A and D are $(l_{ij} - x_{ij})$ and $(x_{ij} - u_{ij})$, respectively.

To describe the algorithm for handling these situations, we need to determine how to form the residual network $G(x)$ for a flow x violating its lower and upper bounds. We consider each arc (i, j) in A one by one and add arcs to the residual network $G(x)$ in the following manner:

1. $l_{ij} \leq x_{ij} \leq u_{ij}$. If $x_{ij} < u_{ij}$, we add the arc (i, j) with a residual capacity $u_{ij} - x_{ij}$ and with a cost c_{ij} . If $x_{ij} > l_{ij}$, we add the arc (j, i) with a residual capacity $x_{ij} - l_{ij}$ and with a cost $-c_{ij}$. We call these arcs *feasible arcs*.
2. $x_{ij} < l_{ij}$. In this case we add the arc (i, j) with a residual capacity $(l_{ij} - x_{ij})$ and with a cost c_{ij} . We refer to this arc as a *lower-infeasible arc*.
3. $x_{ij} > u_{ij}$. In this case we add the arc (j, i) with a residual capacity $(x_{ij} - u_{ij})$ and with a cost $-c_{ij}$. We refer to this arc as an *upper-infeasible arc*.

We next define the kilter numbers of arcs in the residual network. For feasible arcs in the residual network, we define their kilter numbers using (9.23). We define the kilter number k_{ij} of a lower-infeasible or an upper-infeasible arc (i, j) as the change in its residual capacity required to restore its feasibility as well as its optimality. For instance, for a lower-infeasible arc (i, j) (1) if $c_{ij}^\pi \geq 0$, then $k_{ij} = (l_{ij} - x_{ij})$; and (2) if $c_{ij}^\pi < 0$, then $k_{ij} = (u_{ij} - x_{ij})$. Note that

1. Lower-infeasible and upper-infeasible arcs have positive kilter numbers.
2. Sending additional flow on lower-infeasible and upper-infeasible arcs in the residual network decreases their kilter numbers.

The out-of-kilter algorithm for this case is same as that for the earlier case. The algorithmic description given in Figure 9.15 applies to this case as well except that at the beginning of the algorithm we need not establish a feasible flow in the network. We can initiate the algorithm with $x = 0$ as the starting flow. We leave

the justification of the out-of-kilter algorithm for this case as an exercise to the reader (see Exercise 9.26).

9.10 RELAXATION ALGORITHM

All the minimum cost flow algorithms we have discussed so far—the cycle-canceling algorithm, the successive shortest path algorithm, the primal–dual algorithm, and the out-of-kilter algorithm—are classical in the sense that researchers developed them in the 1950s and 1960s as network flow area was emerging as an independent field of scientific investigation. These algorithms have several common features: (1) they repeatedly apply shortest path algorithms, (2) they run in pseudopolynomial time, and (3) their empirical running times have proven to be inferior to those of the network simplex algorithm tailored for the minimum cost flow problem (we discuss this algorithm in Chapter 11). The relaxation algorithm we examine in this section is a more recent vintage minimum cost flow algorithm; it is competitive or better than the network simplex algorithm for some classes of networks. Interestingly, the relaxation algorithm is also a variation of the successive shortest path algorithm. Even though the algorithm has proven to be efficient in practice for many classes of problems, its worst-case running time is much poorer than that of every minimum cost flow algorithm discussed in this chapter.

The relaxation algorithm uses ideas from *Lagrangian relaxation*, a well-known technique used for solving integer programming problems. We discuss the Lagrangian relaxation technique in more detail in Chapter 16. In the Lagrangian relaxation technique, we identify a set of constraints to be relaxed, multiply each such constraint by a scalar, and subtract the product from the objective function. The relaxation algorithm relaxes the mass balance constraints of the nodes, multiplying the mass balance constraint for node i by an (unrestricted) variable $\pi(i)$ (called, as usual, a node potential) and subtracts the resulting product from the objective function. These operations yield the following relaxed problem:

$$w(\pi) = \underset{x}{\text{minimize}} \left[\sum_{(i,j) \in A} c_{ij}x_{ij} + \sum_{i \in N} \pi(i) \left\{ - \sum_{\{j: (i,j) \in A\}} x_{ij} + \sum_{\{j: (j,i) \in A\}} x_{ji} + b(i) \right\} \right] \quad (9.25a)$$

subject to

$$0 \leq x_{ij} \leq u_{ij} \quad \text{for all } (i, j) \in A. \quad (9.25b)$$

For a specific value of the vector π of node potentials, we refer to the relaxed problem as $LR(\pi)$ and denote its objective function value by $w(\pi)$. Note that the optimal solution of $LR(\pi)$ is a pseudoflow for the minimum cost flow problem since it might violate the mass balance constraints. We can restate the objective function of $LR(\pi)$ in the following equivalent way:

$$w(\pi) = \underset{x}{\text{minimize}} \left[\sum_{(i,j) \in A} c_{ij}x_{ij} + \sum_{i \in N} \pi(i)e(i) \right]. \quad (9.26)$$

In this expression, as in our earlier discussion, $e(i)$ denotes the imbalance of node i . Let us restate the objective function (9.25a) of the relaxed problem in another way. Notice that in the second term of (9.25a), each flow variable x_{ij} appears twice: once with a coefficient of $-\pi(i)$ and the second time with a coefficient of $\pi(j)$. Therefore, we can write (9.25a) as follows:

$$w(\pi) = \underset{x}{\text{minimize}} \left[\sum_{(i,j) \in A} (c_{ij} - \pi(i) + \pi(j))x_{ij} + \sum_{i \in N} \pi(i)b(i) \right],$$

or, equivalently,

$$w(\pi) = \underset{x}{\text{minimize}} \left[\sum_{(i,j) \in A} c_{ij}^{\pi} x_{ij} + \sum_{i \in N} \pi(i)b(i) \right]. \quad (9.27)$$

In the subsequent discussion, we refer to the objective function of $\text{LR}(\pi)$ as (9.26) or (9.27), whichever is more convenient. For a given vector π of node potentials, it is very easy to obtain an optimal solution x of $\text{LR}(\pi)$: In light of the formulation (9.27) of the objective function, (1) if $c_{ij}^{\pi} > 0$, we set $x_{ij} = 0$; (2) if $c_{ij}^{\pi} < 0$, we set $x_{ij} = u_{ij}$; and (3) if $c_{ij}^{\pi} = 0$, we can set x_{ij} to any value between 0 and u_{ij} . The resulting solution is a pseudoflow for the minimum cost flow problem and satisfies the reduced cost optimality conditions. We have therefore established the following result.

Property 9.15. *If a pseudoflow x of the minimum cost flow problem satisfies the reduced cost optimality conditions for some π , then x is an optimal solution of $\text{LR}(\pi)$.*

Let z^* denote the optimal objective function value of the minimum cost flow problem. As shown by the next lemma, the value z^* is intimately related to the optimal objective value $w(\pi)$ of the relaxed problem $\text{LR}(\pi)$.

Lemma 9.16

- (a) *For any node potentials π , $w(\pi) \leq z^*$.*
- (b) *For some choice of node potentials π^* , $w(\pi^*) = z^*$.*

Proof. Let x^* be an optimal solution of the minimum cost flow problem with objective function value z^* . Clearly, for any vector π of node potentials, x^* is a feasible solution of $\text{LR}(\pi)$ and its objective function value in $\text{LR}(\pi)$ is also z^* . Therefore, the minimum objective function value of $\text{LR}(\pi)$ will be less than or equal to z^* . We have thus established the first part of the lemma.

To prove the second part, let π^* be a vector of node potentials that together with x^* satisfies the complementary slackness optimality conditions (9.8). Property 9.15 implies that x^* is an optimal solution of $\text{LR}(\pi^*)$ and $w(\pi^*) = cx^* = z^*$. This conclusion completes the proof of the lemma. ♦

Notice the similarity between this result and the weak duality theorem (i.e., Theorem 9.5) for the minimum cost flow problem that we have stated earlier in this chapter. The similarity is more than incidental, since we can view the Lagrangian relaxation solution strategy as a dual linear programming approach that combines

some key features of both the primal and dual linear programs. Moreover, we can view the dual linear program itself as being generated by applying Lagrangian relaxation.

The relaxation algorithm always maintains a vector of node potentials π and a pseudoflow x that is an optimal solution of $\text{LR}(\pi)$. In other words, the pair (x, π) satisfies the reduced cost optimality conditions. The algorithm repeatedly performs one of the following two operations:

1. Keeping π unchanged, it modifies x to x' so that x' is also an optimal solution of $\text{LR}(\pi)$ and the excess of at least one node decreases.
2. It modifies π to π' and x to x' so that x' is an optimal solution of $\text{LR}(\pi')$ and $w(\pi') > w(\pi)$.

If the algorithm can perform either of the two operations, it gives priority to the second operation. Consequently, the primary objective in the relaxation algorithm is to increase $w(\pi)$ and the secondary objective is to reduce the infeasibility of the pseudoflow x while keeping $w(\pi)$ unchanged. We point out that the excesses at the nodes might increase when the algorithm performs the second operation. As we show at the end of this section, these two operations are sufficient to guarantee finite convergence of the algorithm. For a fixed value of $w(\pi)$, the algorithm consistently reduces the excesses of the nodes by at least one unit, and from Lemma 9.16 the number of increases in $w(\pi)$, each of which is at least 1 unit, is finite.

We now describe the relaxation algorithm in more detail. The algorithm performs major iterations and, within a major iteration, it performs several minor iterations. Within a major iteration, the algorithm selects an excess node s and grows a tree rooted at node s so that every tree node has a nonnegative imbalance and every tree arc has zero reduced cost. Each minor iteration adds an additional node to the tree. A major iteration ends when the algorithm performs either an augmentation or increases $w(\pi)$.

Let S denote the set of nodes spanned by the tree at some stage and let $\bar{S} = N - S$. The set S defines a cut which we denote by $[S, \bar{S}]$. As in earlier chapters, we let (S, \bar{S}) denote the set of forward arcs in the cut and (\bar{S}, S) the set of backward arcs [all in $G(x)$]. The algorithm maintains two variables $e(S)$ and $r(\pi, S)$, defined as follows:

$$e(S) = \sum_{i \in S} e(i),$$

$$r(\pi, S) = \sum_{(i,j) \in (S, \bar{S}) \text{ and } c_{ij}^x = 0} r_{ij}.$$

Given the set S , the algorithm first checks the condition $e(S) > r(\pi, S)$. If the current solution satisfies this condition, the algorithm can increase $w(\pi)$ in the following manner. [We illustrate this method using the example shown in Figure 9.17(a).] The algorithm first increases the flow on zero reduced cost arcs in (S, \bar{S}) so that they become saturated (i.e., drop out of the residual network). The flow change does not alter the value of $w(\pi)$ because the change takes place on arcs with zero reduced costs. However, the flow change decreases the total imbalance of the

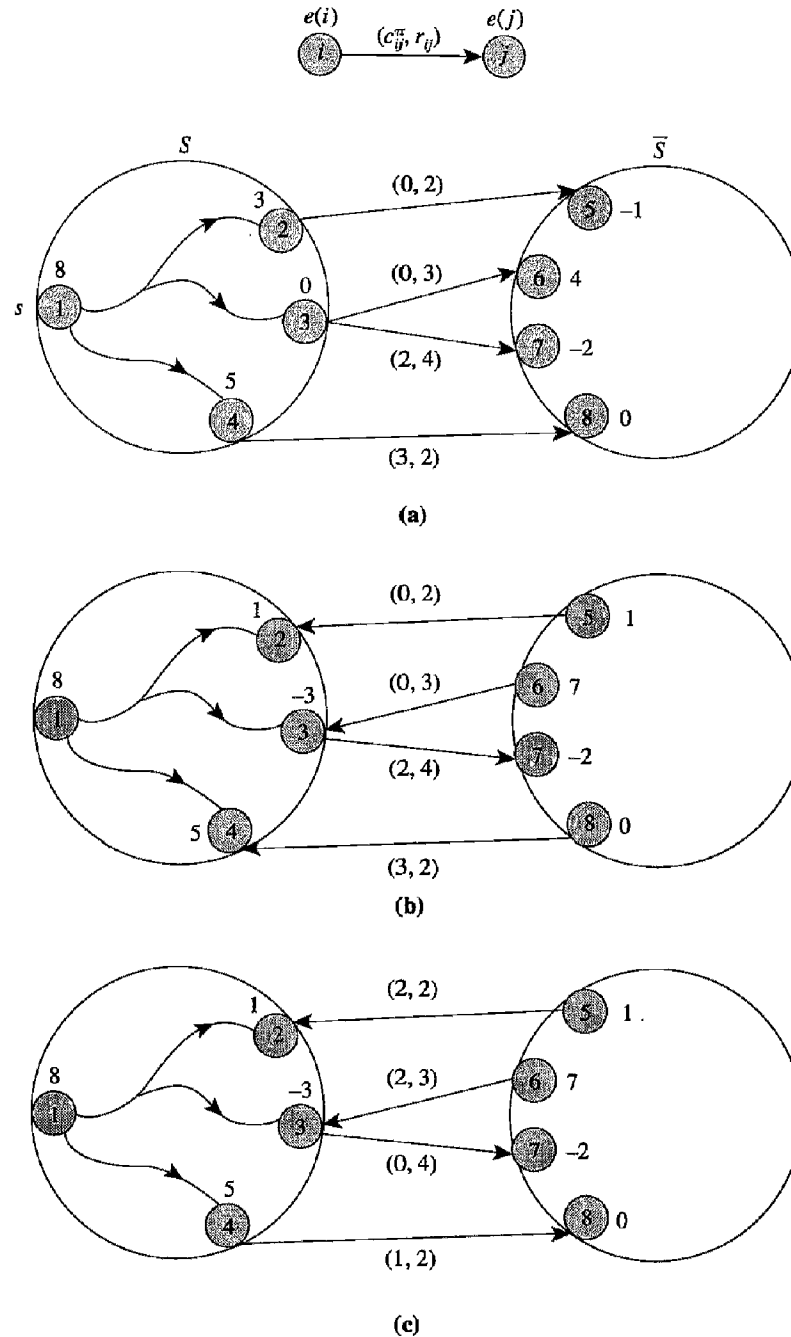


Figure 9.17 Illustrating the relaxation algorithm: (a) solution at some stage; (b) solution after modifying the flow; (c) solution after modifying the potentials.

nodes by the amount $r(\pi, S)$; but since $e(S) > r(\pi, S)$, the remaining imbalance $e(S) - r(\pi, S)$ is still positive [see Figure 9.17(b)].

At this point all the arcs in (S, \bar{S}) have (strictly) positive reduced cost. The algorithm next computes the minimum reduced cost of an arc in (S, \bar{S}) , say α , and increases the potential of every node $i \in S$ by $\alpha > 0$ units [see Figure 9.17(c)]. The

formulation (9.26) of the Lagrangian relaxation objective function implies that this updating of the node potentials does not change its first term but increases the second term by $(e(S) - r(\pi, S))\alpha$ units. Therefore, this operation increases $w(\pi)$ by $(e(S) - r(\pi, S))\alpha$ units, which is strictly positive. Increasing the potentials of nodes in S by α decreases the reduced costs of all the arcs in (S, \bar{S}) by α units, increases the reduced costs of all arcs in (\bar{S}, S) by α units, and does not change the remaining reduced costs. Although increasing the reduced costs does not change the reduced cost optimality conditions, decreasing the reduced costs might. Notice, however, that before we change the node potentials, $c_{ij}^{\pi} \geq \alpha$ for all $(i, j) \in (S, \bar{S})$; therefore, after the change, $c_{ij}^{\pi} \geq 0$, so the algorithm preserves the optimality conditions. This completes one major iteration.

We next study situations in which $e(S) \leq r(\pi, S)$. Since $r(\pi, S) \geq e(S) > 0$, at least one arc $(i, j) \in (S, \bar{S})$ must have a zero reduced cost. If $e(j) \geq 0$, the algorithm adds node j to S , completes one minor iteration, and repeats this process. If $e(j) < 0$, the algorithm augments the maximum possible flow along the tree path from node s to node j . Notice that since we augment flow along zero residual cost arcs, we do not change the objective function value of $LR(\pi)$. The augmentation reduces the total excess of the nodes and completes one major iteration of the algorithm.

Figures 9.18 and 9.19 give a formal description of the relaxation algorithm.

It is easy to see that the algorithm terminates with a minimum cost flow. The algorithm terminates when all of the node imbalances have become zero (i.e., the solution is a flow). Because the algorithm maintains the reduced cost optimality conditions at every iteration, the terminal solution is a minimum cost flow.

We now prove that for problems with integral data, the algorithm terminates in a finite number of iterations. Since each minor iteration adds a node to the set S , within n minor iterations the algorithm either executes adjust-flow or executes adjust-potentials. Each call of the procedure adjust-flow decreases the excess of at least one node by at least 1 unit; therefore, the algorithm can perform a finite number of executions of the adjust-flow procedure within two consecutive calls of the adjust-potential procedure. To bound the executions of the adjust-potential procedure, we notice that (1) initially, $w(\pi) = 0$; (2) each call of this procedure strictly increases

```

algorithm relaxation;
begin
   $x := 0$  and  $\pi := 0$ ;
  while the network contains a node  $s$  with  $e(s) > 0$  do
    begin
       $S := \{s\}$ ;
      if  $e(S) > r(\pi, S)$  then adjust-potential;
      repeat
        select an arc  $(i, j) \in (S, \bar{S})$  in the residual network with  $c_{ij}^{\pi} = 0$ ;
        if  $e(j) \geq 0$  then set  $\text{pred}(j) := i$  and add node  $j$  to  $S$ ;
      until  $e(j) < 0$  or  $e(S) > r(\pi, S)$ ;
      if  $e(S) > r(\pi, S)$  then adjust-potential
      else adjust-flow;
    end;
  end;

```

Figure 9.18 Relaxation algorithm.

```

procedure adjust-potential;
begin
  for every arc  $(i, j) \in (S, \bar{S})$  with  $c_{ij}^r = 0$  do send  $r_{ij}$  units of flow on the arc  $(i, j)$ ;
  compute  $\alpha := \min\{c_{ij}^r : (i, j) \in (S, \bar{S}) \text{ and } r_{ij} > 0\}$ ;
  for every node  $i \in S$  do  $\pi(i) := \pi(i) + \alpha$ ;
end;

(a)

procedure adjust-flow;
begin
  trace the predecessor indices to identify the directed path  $P$  from node  $s$  to node  $j$ ;
   $\delta := \min[e(s) - e(j), \min\{r_{ij} : (i, j) \in P\}]$ ;
  augment  $\delta$  units of flow along  $P$ , update imbalances and residual capacities;
end;

(b)

```

Figure 9.19 Procedures of the relaxation algorithm.

$w(\pi)$ by at least 1 unit; and (3) the maximum possible value of $w(\pi)$ is mCU . The preceding arguments establish that the algorithm performs finite number of iterations. In Exercise 9.27 we ask the reader to obtain a worst-case bound on the total number of iterations; this time bound is much worse than those of the other minimum cost flow algorithms discussed in earlier sections.

Notice that the relaxation algorithm is a type of shortest augmenting path algorithm; indeed, it bears some resemblance to the successive shortest path algorithm that we considered in Section 9.7. Since the reduced cost of every arc in the residual network is nonnegative, and since every arc in the tree connecting the nodes in S has a zero reduced cost, the path P that we find in the adjust-flow procedure of the relaxation algorithm is a shortest path in the residual network. Therefore, the sequence of flow adjustments that the algorithm makes is a set of flow augmentations along shortest augmenting paths. The relaxation algorithm differs from the successive shortest augmenting path algorithm, however, because it uses “intermediate” information to make changes to the node potentials as it fans out and constructs the tree containing the nodes S . This use of intermediate information might explain why the relaxation algorithm has performed much better empirically than the successive shortest path algorithm.

9.11 SENSITIVITY ANALYSIS

The purpose of sensitivity analysis is to determine changes in the optimal solution of a minimum cost flow problem resulting from changes in the data (supply/demand vector or the capacity or cost of any arc). There are two different ways of performing sensitivity analysis: (1) using combinatorial methods, and (2) using simplex-based methods from linear programming. Each method has its advantages. For example, although combinatorial methods obtain better worst-case time bounds for performing sensitivity analysis, simplex-based methods might be more efficient in practice. In this section we describe sensitivity analysis using combinatorial methods; in Section

11.10 we consider a simplex-based approach. For simplicity, we limit our discussion to a unit change of only a particular type. In a sense, however, this discussion is quite general: It is possible to reduce more complex changes to a sequence of the simple changes we consider. We show that sensitivity analysis for the minimum cost flow problem essentially reduces to applying shortest path or maximum flow algorithms.

Let x^* denote an optimal solution of a minimum cost flow problem. Let π be the corresponding node potentials and $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$ denote the reduced costs. Further, let $d(k, l)$ denote the shortest distance from node k to node l in the residual network with respect to the original arc lengths c_{ij} . Since for any directed path P from node k to node l , $\sum_{(i,j) \in P} c_{ij}^\pi = \sum_{(i,j) \in P} c_{ij} - \pi(k) + \pi(l)$, $d(k, l)$ equals the shortest distance from node k to node l with respect to the arc lengths c_{ij}^π plus $[\pi(k) - \pi(l)]$. At optimality, the reduced costs c_{ij}^π of all arcs in the residual network are nonnegative. Therefore, we can compute $d(k, l)$ for all pairs of nodes k and l by solving n single-source shortest path problems with nonnegative arc lengths.

Supply/Demand Sensitivity Analysis

We first study changes in the supply/demand vector. Suppose that the supply/demand of a node k becomes $b(k) + 1$ and the supply/demand of another node l becomes $b(l) - 1$. [Recall from Section 9.1 that feasibility of the minimum cost flow problem dictates that $\sum_{i \in N} b(i) = 0$; therefore, we must change the supply/demand values of two nodes by equal magnitudes, and must increase one value and decrease the other.] The vector x^* is a pseudoflow for the modified problem; moreover, this vector satisfies the reduced cost optimality conditions. Augmenting 1 unit of flow from node k to node l along the shortest path in the residual network $G(x^*)$ converts this pseudoflow into a flow. This augmentation changes the objective function value by $d(k, l)$ units. Lemma 9.12 implies that this flow is optimal for the modified minimum cost flow problem. We point out that the residual network $G(x^*)$ might not contain any directed path from node k to node l , in which case the modified minimum cost flow problem is infeasible.

Arc Capacity Sensitivity Analysis

We next consider a change in an arc capacity. Suppose that the capacity of an arc (p, q) increases by 1 unit. The flow x^* is feasible for the modified problem. In addition, if $c_{pq}^\pi \geq 0$, it satisfies the reduced cost optimality conditions; therefore, it is an optimal flow for the modified problem. If $c_{pq}^\pi < 0$, the optimality conditions dictate that the flow on the arc must equal its capacity. We satisfy this requirement by increasing the flow on the arc (p, q) by 1 unit, which produces a pseudoflow with an excess of 1 unit at node q and a deficit of 1 unit at node p . We convert the pseudoflow into a flow by augmenting 1 unit of flow from node q to node p along the shortest path in the residual network $G(x^*)$, which changes the objective function value by an amount $c_{pq} + d(q, p)$. This flow is optimal from our observations concerning supply/demand sensitivity analysis.

When the capacity of the arc (p, q) decreases by 1 unit and the flow on the arc is strictly less than its capacity, x^* remains feasible, and therefore optimal, for the modified problem. However, if the flow on the arc is at its capacity, we decrease

the flow by 1 unit and augment 1 unit of flow from node p to node q along the shortest path in the residual network. This augmentation changes the objective function value by an amount $-c_{pq} + d(p, q)$. Observed that the residual network $G(x^*)$ might not contain any directed path from node p to node q , indicating the infeasibility of the modified problem.

Cost Sensitivity Analysis

Finally, we discuss changes in arc costs, which we assume are integral. We discuss the case when the cost of an arc (p, q) increases by 1 unit; the case when the cost of an arc decreases is left as an exercise to the reader (see Exercise 9.50). This change increases the reduced cost of arc (p, q) by 1 unit as well. If $c_{pq}^\pi < 0$ before the change, then after the change, the modified reduced cost is nonpositive. Similarly, if $c_{pq}^\pi > 0$ before the change, the modified reduced cost is nonnegative after the change. In both cases we preserve the optimality conditions. However, if $c_{pq}^\pi = 0$ before the change and $x_{pq} > 0$, then after the change the modified reduced cost is positive and the solution violates the reduced-cost optimality conditions. To restore the optimality conditions of the arc, we must either reduce the flow on arc (p, q) to zero or change the potentials so that the reduced cost of arc (p, q) becomes zero.

We first try to reroute the flow x_{pq}^* from node p to node q without violating any of the optimality conditions. We do so by solving a maximum flow problem defined as follows: (1) set the flow on the arc (p, q) to zero, thus creating an excess of x_{pq}^* at node p and a deficit of x_{pq}^* at node q ; (2) designate node p as the source node and node q as the sink node; and (3) send a maximum of x_{pq}^* units from the source to the sink. We permit the maximum flow algorithm, however, to change flows only on arcs with zero reduced costs since otherwise it would generate a solution that might violate (9.8). Let v° denote the flow sent from node p to node q and x° denote the resulting arc flow. If $v^\circ = x_{pq}^*$, then x° denotes a minimum cost flow of the modified problem. In this case the optimal objective function values of the original and modified problems are the same.

On the other hand, if $v^\circ < x_{pq}^*$, the maximum flow algorithm yields an s - t cut $[S, \bar{S}]$ with the properties that $p \in S$, $q \in \bar{S}$, and every forward arc in the cut with zero reduced cost has flow equal to its capacity and every backward arc in the cut with zero reduced cost has zero flow. We then decrease the node potential of every node in \bar{S} by 1 unit. It is easy to verify by case analysis that this change in node potentials maintains the complementary slackness optimality conditions and, furthermore, decreases the reduced cost of arc (p, q) to zero. Consequently, we can set the flow on arc (p, q) equal to $x_{pq}^* - v^\circ$ and obtain a feasible minimum cost flow. In this case the objective function value of the modified problem is $x_{pq}^* - v^\circ$ units more than that of the original problem.

9.12 SUMMARY

The minimum cost flow problem is the central object of study in this book. In this chapter we began our study of this important class of problems by showing how minimum cost flow problems arise in several application settings and by considering

Algorithm	Number of iterations	Features
Cycle-canceling algorithm	$O(mCU)$	<ol style="list-style-type: none"> 1. Maintains a feasible flow x at every iteration and augments flows along negative cycles in $G(x)$. 2. At each iteration, solves a shortest path problem with arbitrary arc lengths to identify a negative cycle. 3. Very flexible: some rules for selecting negative cycles leads to polynomial-time algorithms.
Successive shortest path algorithm	$O(nU)$	<ol style="list-style-type: none"> 1. Maintains a pseudoflow x satisfying the optimality conditions and augments flow along shortest paths from excess nodes to deficit nodes in $G(x)$. 2. At each iteration, solves a shortest path problem with non-negative arc lengths. 3. Very flexible: by selecting augmentations carefully, we can obtain several polynomial-time algorithms.
Primal-dual algorithm	$O(\min\{nU, nC\})$	<ol style="list-style-type: none"> 1. Maintains a pseudoflow x satisfying the optimality conditions. Solves a shortest path problem to update node potentials and attempts to reduce primal infeasibility by the maximum amount by solving a maximum flow problem. 2. At each iteration, solves both a shortest path problem with nonnegative arc lengths and a maximum flow problem. 3. Closely related to the successive shortest path algorithm: instead of sending flow along one shortest path, sends flow along all shortest paths.
Out-of-kilter algorithm	$O(nU)$	<ol style="list-style-type: none"> 1. Maintains a feasible flow x at each iteration and attempts to satisfy the optimality conditions by augmenting flows along shortest paths. 2. At each iteration, solves a shortest path problem with non-negative arc lengths. 3. Can be generalized to solve situations in which the flow x maintained by the algorithm might not satisfy the flow bounds on the arcs.
Relaxation algorithm	See Exercise 9.27	<ol style="list-style-type: none"> 1. Somewhat different from other minimum cost flow algorithms. 2. Maintains a pseudoflow x satisfying the optimality conditions and modifies arc flows and node potentials so that a Lagrangian objective function does not decrease and occasionally increases. 3. With the incorporation of some heuristics, the algorithm is very efficient in practice and yields the fastest available algorithm for some classes of minimum cost flow problems.

Figure 9.20 Summary of pseudopolynomial-time algorithms for the minimum cost flow problem.

the simplest pseudopolynomial-time algorithms for solving these problems. These pseudopolynomial-time algorithms include classical algorithms that are important because of both their historical significance and because they provide the essential building blocks and core ideas used in more efficient algorithms. Our algorithmic development relies heavily upon optimality conditions for the minimum cost flow problem that we developed and proved in the following equivalent frameworks: negative cycle optimality conditions, reduced cost optimality conditions, and complementary slackness optimality conditions. The negative cycle optimality conditions state that a feasible flow x is an optimal flow if and only if the residual network $G(x)$ contains no negative cycle. The reduced cost optimality conditions state that a feasible flow x is an optimal flow if and only if the reduced cost of each arc in the residual network is nonnegative. The complementary slackness optimality conditions are adaptations of the linear programming optimality conditions for network flows. As part of this general discussion in this chapter, we also examined minimum cost flow duality.

We developed several minimum cost flow algorithms: the cycle-canceling, successive shortest path, primal–dual, out-of-kilter, and relaxation algorithms. These algorithms represent a good spectrum of approaches for solving the same problem: Some of these algorithms maintain primal feasible solutions and strive toward optimality; others maintain primal infeasible solutions that satisfy the optimality conditions and strive toward feasibility. These algorithms have some commonalities as well—they all repeatedly solve shortest path problems. In fact, in Exercises 9.57 and 9.58 we establish a very strong result by showing that the cycle-canceling, successive shortest path, primal–dual, and out-of-kilter algorithms are all equivalent in the sense that if initialized properly, they perform the same sequence of augmentations. Figure 9.20 summarizes the basic features of the algorithms discussed in this chapter.

Finally, we discussed sensitivity analysis for the minimum cost flow problem. We showed how to reoptimize the minimum cost flow problem, after we have made unit changes in the supply/demand vector or the arc capacities; by solving a shortest path problem, and how to handle unit changes in the cost vector by solving a maximum flow problem. Needless to say, these reoptimization procedures are substantially faster than solving the problem afresh if the changes in the problem data are sufficiently small.

REFERENCE NOTES

In this chapter and in these reference notes we focus on pseudopolynomial-time nonsimplex algorithms for solving minimum cost flow problems. In Chapter 10 we provide references for polynomial-time minimum cost flow algorithms, and in Chapter 11 we give references for simplex-based algorithms.

Ford and Fulkerson [1957] developed the primal–dual algorithms for the capacitated transportation problem; Ford and Fulkerson [1962] later generalized this approach for solving the minimum cost flow problem. Jewell [1958], Iri [1960], and Busaker and Gowen [1961] independently developed the successive shortest path algorithm. These researchers showed how to solve the minimum cost flow problem as a sequence of shortest path problems with arbitrary arc lengths. Tomizawa [1972]

and Edmonds and Karp [1972] independently observed that if the computations use node potentials, it is possible to implement these algorithms so that the shortest path problems have nonnegative arc lengths.

Minty [1960] and Fulkerson [1961b] independently developed the out-of-kilter algorithm. Aashtiani and Magnanti [1976] have described an efficient implementation of this algorithm. The description of the out-of-kilter algorithm presented in Section 9.9 differs substantially from the development found in other textbooks. Our description is substantially shorter and simpler because it avoids tedious case analyses. Moreover, our description explicitly highlights the use of Dijkstra's algorithm; because other descriptions do not focus on the shortest path computations, they find an accurate worst-case analysis of the algorithm much more difficult to conduct.

The cycle-canceling algorithm is credited to Klein [1967]. Three special implementations of the cycle-canceling algorithms run in polynomial time: the first, due to Barahona and Tardos [1989] (which, in turn, modifies an algorithm by Weintraub [1974]), augments flow along (negative) cycles with the maximum possible improvement; the second, due to Goldberg and Tarjan [1988], augments flow along minimum mean cost (negative) cycles; and the third, due to Wallacher and Zimmerman [1991], augments flow along minimum ratio cycles.

Zadeh [1973a,1973b] described families of minimum cost flow problems on which each of several algorithms—the cycle-canceling algorithm, successive shortest path algorithm, primal–dual algorithm, and out-of-kilter algorithm—perform an exponential number of iterations. The fact that the same families of networks are bad for many network algorithms suggests an interrelationship among the algorithms. The insightful paper by Zadeh [1979] points out that each of the algorithms we have just mentioned are indeed equivalent in the sense that they perform the same sequence of augmentations, which they obtained through shortest path computations, provided that we initialize them properly and break ties using the same rule.

Bertsekas and Tseng [1988b] developed the relaxation algorithm and conducted extensive computational investigations of it. A FORTRAN code of the relaxation algorithm appears in Bertsekas and Tseng [1988a]. Their study and those conducted by Grigoriadis [1986] and Kennington and Wang [1990] indicate that the relaxation algorithm and the network simplex algorithm (described in Chapter 11) are the two fastest available algorithms for solving the minimum cost flow problem in practice. When the supplies/demands at nodes are relatively small, the successive shortest path algorithm is the fastest algorithm. Previous computational studies conducted by Glover, Karney, and Klingman [1974] and Bradley, Brown, and Graves [1977] have indicated that the network simplex algorithm is consistently superior to the primal–dual and out-of-kilter algorithms. Most of these computational testings have been done on random network flow problems generated by the well-known computer program NETGEN, suggested by Klingman, Napier, and Stutz [1974].

The applications of the minimum cost flow problem that we discussed Section 9.2 have been adapted from the following papers:

1. Distribution problems (Glover and Klingman [1976])
2. Reconstructing the left ventricle from x-ray projections (Slump and Gerbrands [1982])
3. Racial balancing of schools (Belford and Ratliff [1972])

4. Optimal loading of a hopping airplane (Gupta [1985] and Lawania [1990])
5. Scheduling with deferral costs (Lawler [1964])
6. Linear programming with consecutive 1's in columns (Veinott and Wagner [1962])

Elsewhere in this book we describe other applications of the minimum cost flow problem. These applications include (1) leveling mountainous terrain (Application 1.4, Farley [1980]), (2) the forest scheduling problem (Exercise 1.10), (3) the entrepreneur's problem (Exercise 9.1, Prager [1957]), (4) vehicle fleet planning (Exercise 9.2), (5) optimal storage policy for libraries (Exercise 9.3, Evans [1984]), (6) zoned warehousing (Exercise 9.4, Evans [1984]), (7) allocation of contractors to public works (Exercise 9.5, Cheshire, McKinnon, and Williams [1984]), (8) phasing out capital equipment (Exercise 9.6, Daniel [1973]), (9) the terminal assignment problem (Exercise 9.7, Esau and Williams [1966]), (10) linear programs with consecutive or circular 1's in rows (Exercises 9.8 and 9.9, Bartholdi, Orlin, and Ratliff [1980]), (11) capacitated maximum spanning trees (Exercise 9.54, Garey and Johnson [1979]), (12) fractional b -matching (Exercise 9.55), (13) the nurse scheduling problem (Exercise 11.1), (14) the caterer problem (Exercise 11.2, Jacobs [1954]), (15) project assignment (Exercise 11.3), (16) passenger routing (Exercise 11.4), (17) allocating receivers to transmitters (Exercise 11.5, Dantzig [1962]), (18) faculty-course assignment (Exercise 11.6, Mulvey [1979]), (19) optimal rounding of a matrix (Exercise 11.7, Bacharach [1966], Cox and Ernst [1982]), (20) automatic karyotyping of chromosomes (Application 19.8, Tso, Kleinschmidt, Mitterreiter, and Graham [1991]), (21) just-in-time scheduling (Application 19.10, Elmaghraby [1978], Levner and Nemirovsky [1991]), (22) time-cost trade-off in project management (Application 19.11, Fulkerson [1961a] and Kelly [1961]), (23) models for building evacuation (Application 19.13, Chalmet, Francis and Saunders [1982]), (24) the directed Chinese postman problem (Application 19.14, Edmonds and Johnson [1973]), (25) warehouse layout (Application 19.17, Francis and White [1976]), (26) rectilinear distance facility location (Application 19.18, Cabot, Francis, and Stary [1970]), (27) dynamic lot sizing (Application 19.19, Zangwill [1969]), (28) multistage production-inventory planning (Application 19.23, Evans [1977]), (29) mold allocation (Application 19.24, Love and Vemuganti [1978]), (30) a parking model (Exercise 19.17, Dirickx and Jennergren [1975]), (31) the network interdiction problem (Exercise 19.18, Fulkerson and Harding [1977]), (32) truck scheduling (Exercises 19.19 and 19.20, Gavish and Schweitzer [1974]), and (33) optimal deployment of firefighting companies (Exercise 19.21, Denardo, Rothblum, and Swersey [1988]).

The applications of the minimum cost flow problems are so vast that we have not been able to describe many other applications in this book. The following list provides a set of references to some other applications: (1) warehousing and distribution of a seasonal product (Jewell [1957]), (2) economic distribution of coal supplies in the gas industry (Berrisford [1960]), (3) upsets in round-robin tournaments (Fulkerson [1965]), (4) optimal container inventory and routing (Horn [1971]), (5) distribution of empty rail containers (White [1972]), (6) optimal defense of a network (Picard and Ratliff [1973]), (7) telephone operator scheduling (Segal [1974]), (8) multifacility minimax location problem with rectilinear distances (Dearing and Francis [1974]), (9) cash management problems (Srinivasan [1974]), (10) multiproduct mul-

tifacility production-inventory planning (Dorsey, Hodgson, and Ratliff [1975]), (11) “hub” and “wheel” scheduling problems (Arisawa and Elmaghraby [1977]), (12) the warehouse leasing problem (Lowe, Francis, and Reinhardt [1979]), (13) multiattribute marketing models (Srinivasan [1979]), (14) material handling systems (Maxwell and Wilson [1981]), (15) microdata file merging (Barr and Turner [1981]), (16) determining service districts (Larson and Odoni [1981]), (17) control of forest fires (Kourtz [1984]), (18) allocating blood to hospitals from a central blood bank (Sapountzis [1984]), (19) market equilibrium problems (Dafetmos and Nagurney [1984]), (20) automatic chromosome classifications (Tso [1986]), (21) the city traffic congestion problem (Zawack and Thompson [1987]), (22) satellite scheduling (Servi [1989]), and (23) determining k disjoint cuts in a network (Wagner [1990]).

EXERCISES

- 9.1. Entrepreneur’s problem** (Prager [1957]). An entrepreneur faces the following problem. In each of T periods, he can buy, sell, or hold for later sale some commodity, subject to the following constraints. In each period i he can buy at most α_i units of the commodity, can holdover at most β_i units of the commodity for the next period, and must sell at least γ_i units (perhaps due to prior agreements). The entrepreneur cannot sell the commodity in the same period in which he buys it. Assuming that p_i , w_i , and s_i denote the purchase cost, inventory carrying cost, and selling price per unit in period i , what buy–sell policy should the entrepreneur adopt to maximize total profit in the T periods? Formulate this problem as a minimum cost flow problem for $T = 4$.
- 9.2. Vehicle fleet planning.** The Millersburg Supply Company uses a large fleet of vehicles which it leases from manufacturers. The company has forecast the following pattern of vehicle requirements for the next 6 months:

Month	Jan.	Feb.	Mar.	Apr.	May	June
Vehicles required	430	410	440	390	425	450

Millersburg can lease vehicles from several manufacturers at various costs and for various lengths of time. Three of the plans appear to be the best available: a 3-month lease for \$1700; a 4-month lease for \$2200; and a 5-month lease for \$2600. The company can undertake a lease beginning in any month. On January 1 the company has 200 cars on lease, all of which go off lease at the end of February. Formulate the problem of determining the most economical leasing policy as a minimum cost flow problem. (*Hint:* Observe that the linear (integer) programming formulation of this problem has consecutive 1’s in each column. Then use the result in Application 9.6.)

- 9.3. Optimal storage policy for libraries** (Evans [1984]). A library facing insufficient primary storage space for its collection is considering the possibility of using secondary facilities, such as closed stacks or remote locations, to store portions of its collection. These options are preferred to an expensive expansion of primary storage. Each secondary storage facility has limited capacity and a particular access costs for retrieving information. Through appropriate data collection, we can determine the usage rates for the information needs of the users. Let b_j denote the capacity of storage facility j and v_j

denote the access cost per unit item from this facility. In addition, let a_i denote the number of items of a particular class i requiring storage and let u_i denote the expected rate (per unit time) that we will need to retrieve books from this class. Our goal is to store the books in a way that will minimize the expected retrieval cost.

- (a) Show how to formulate the problem of determining an optimal policy as a transportation problem. What is the special structure of this problem? Transportation problems with this structure have become known as *factored transportation problems*.
- (b) Show that the simple rule that repeatedly assigns items with the greatest retrieval rate to the storage facility with lowest access cost specifies an optimal solution of this library storage problem.

9.4. Zoned warehousing (Evans [1984]). In the storage of multiple, say p , items in a zoned warehouse, we need to extract (pick) items in large quantities (perhaps by pallet loads). Suppose that the warehouse is partitioned into q zones, each with a different distance to the shipping area. Let B_j denote the storage capacity of zone j and let d_j denote the average distance from zone j to the shipping area. For each item i , we know (1) the space requirement per unit (r_i), (2) the average order size in some common volume unit (s_i), and (3) the average number of orders per day (f_i). The problem is to determine the quantity of each item to allocate to each zone in order to minimize the average daily handling costs. Assume that the handling cost is linearly proportional to the distance and to the volume moved.

- (a) Formulate this problem as a factored transportation problem (as defined in Exercise 9.3).
- (b) Specify a simple rule that yields an optimal solution of the zoned warehousing problem.

9.5. Allocation of contractors to public works (Cheshire, McKinnon, and Williams [1984]). A large publicly owned corporation has 12 divisions in Great Britain. Each division faces a similar problem. Each year the division subcontracts work to private contractors. The work is of several different types and is done by teams, each of which is capable of doing all types of work. One of these divisions is divided into several districts: the j th district requires r_j teams. The contractors are of two types: experienced and inexperienced. Each contractor i quotes a price c_{ij} to have a team conduct the work in district j . The objective is to allocate the work in the districts to the various contractors, satisfying the following conditions: (1) each district j has r_j assigned teams; (2) the division contracts with contractor i for no more than u_i teams, the maximum number of teams it can supply; and (3) each district has at least one experienced contractor assigned to it. Formulate this problem as a minimum cost flow problem for a division with three districts, and with two experienced and two inexperienced contractors. (*Hint*: Split each district node into two nodes, one of which requires an experienced contractor.)

9.6. Phasing out capital equipment (Daniel [1973]). A shipping company wants to phase out a fleet of (homogeneous) general cargo ships over a period of p years. Its objective is to maximize its cash assets at the end of the p years by considering the possibility of prematurely selling ships and temporarily replacing them by charter ships. The company faces a known nonincreasing demand for ships. Let $d(i)$ denote the demand of ships in year i . Each ship earns a revenue of r_k units in period k . At the beginning of year k , the company can sell any ship that it owns, accruing a cash inflow of s_k dollars. If the company does not own sufficiently many ships to meet its demand, it must hire additional charter ships. Let h_k denote the cost of hiring a ship for the k th year. The shipping company wants to meet its commitments and at the same time maximize the cash assets at the end of the p th year. Formulate this problem as a minimum cost flow problem.

- 9.7. Terminal assignment problem** (Esau and Williams [1966]). Centralized teleprocessing networks often contain many (as many as tens of thousands) relatively unsophisticated geographically dispersed terminals. These terminals need to be connected to a central processor unit (CPU) either by direct lines or through *concentrators*. Each concentrator is connected to the CPU through a high-speed, cost-effective line that is capable of merging data flow streams from different terminals and sending them to the CPU. Suppose that the concentrators are in place and that each concentrator can handle at most K terminals. For each terminal j , let c_{oj} denote the cost of laying down a direct line from the CPU to the terminal and let c_{ij} denote the line construction cost for connecting concentrator i to terminal j . The decision problem is to construct the minimum cost network for connecting the terminals to the CPU. Formulate this problem as a minimum cost flow problem.
- 9.8. Linear programs with consecutive 1's in rows.** In Application 9.6 we considered linear programs with consecutive 1's in each column and showed how to transform them into minimum cost flow problems. In this and the next exercise we study several related linear programming problems and show how we can solve them by solving minimum cost flow problems. In this exercise we study linear programs with consecutive 1's in the rows. Consider the following (integer) linear program with consecutive 1's in the rows:

$$\text{Minimize } c_1x_1 + c_2x_2 + c_3x_3 + c_4x_4$$

subject to

$$\begin{aligned} x_2 + x_3 + x_4 &\geq 20 \\ x_1 + x_2 + x_3 + x_4 &\geq 30 \\ x_2 + x_3 &\geq 15 \\ x_3 + x_4 &\geq 10 \\ x_1, x_2, x_3, x_4 &\geq 0 \text{ and integer.} \end{aligned}$$

Transform this problem to a minimum cost flow problem. (*Hint:* Use the same transformation of variables that we used in Application 4.6.)

- 9.9. Linear programs with circular 1's in rows** (Bartholdi, Orlin, and Ratliff [1980]). In this exercise we consider a generalization of Exercise 9.8 with the 1's in each row arranged consecutively when we view columns in the wraparound fashion (i.e., we consider the first column as next to the last column). A special case of this problem is the telephone operator scheduling problem that we discussed in Application 4.6. In this exercise we focus on the telephone operator scheduling problem; nevertheless, the approach easily extends to any general linear program with circular 1's in the rows. We consider a version of the telephone operator scheduling in which we incur a cost c_i whenever an operator works in the i th shift, and we wish to satisfy the minimum operator requirement for each hour of the day at the least possible cost. We can formulate this “cyclic staff scheduling problem” as the following (integer) linear program.

$$\text{Minimize } \sum_{i=0}^{23} y_i$$

subject to

$$\begin{aligned} y_{i-7} + y_{i-6} + \cdots + y_i &\geq b(i) && \text{for all } i = 7 \text{ to } 23, \\ y_{17+i} + \cdots + y_{23} + y_0 + \cdots + y_i &\geq b(i) && \text{for all } i = 0 \text{ to } 6, \\ y_i &\geq 0 && \text{for all } i = 1 \text{ to } 23. \end{aligned}$$

- (a) For a parameter p , let $\mathcal{P}(p)$ denote the cyclic staff scheduling problem when we impose the additional constraint $\sum_{i=0}^{23} y_i = p$, and let $z(p)$ denote the optimal objective value of this problem. Show how to transform $\mathcal{P}(p)$, for a fixed value of p , into a minimum cost flow problem. (*Hint*: Use the same transformation of variables that we used in Application 4.6 and observe that each row has one $+1$ and one -1 . Then use the result of Theorem 9.9.)
 - (b) Show that $z(p)$ is a (piecewise linear) convex function of p . (*Hint*: Show that if y' is an optimal solution of $\mathcal{P}(p')$ and y'' is an optimal solution of $\mathcal{P}(p'')$, then for any weighting parameter λ , $0 \leq \lambda \leq 1$, the point $\lambda y' + (1 - \lambda)y''$ is a feasible solution of $\mathcal{P}(\lambda p' + (1 - \lambda)p'')$.)
 - (c) In the cyclic staff scheduling problem, we wish to determine a value of p , say p^* , satisfying the property that $z(p^*) \leq z(p)$ for all feasible p . Show how to solve the cyclic staff scheduling problem in polynomial time by performing binary search on the values of p . (*Hint*: For any integer p , show how to determine whether $p \leq p^*$ by solving problems $\mathcal{P}(p)$ and $\mathcal{P}(p + 1)$.)
- 9.10. Racial balancing of schools.** In this exercise we discuss some generalizations of the problem of racial balancing of schools that we described in Application 9.3. Describe how would you modify the formulation to include the following additional restrictions (consider each restriction separately).
- (a) We prohibit the assignment of a student from location i to school j if the travel distance d_{ij} between these location exceeds some specified distance D .
 - (b) We include the distance traveled between location i and school j in the objective function only if d_{ij} is greater than some specified distance D' (e.g., we account for the distance traveled only if a student needs to be bussed).
 - (c) We impose lower and upper bounds on the number of black students from location i who are assigned to school j .
- 9.11.** Show how to transform the equipment replacement problem described in Application 9.6 into a shortest path problem. Give the resulting formulation for $n = 4$.
- 9.12.** This exercise is based on the equipment replacement problem that we discussed in Application 9.6.
- (a) The problem as described allows us to buy and sell the equipment only yearly. How would you model the situation if you could make decisions every half year?
 - (b) How sensitive do you think the optimal solution would be to the length T of planning period? Can you anticipate a situation in which the optimal replacement plan would change drastically if we were to increase the length of the planning period to $T + 1$?
- 9.13.** Justify the minimum cost flow formulation that we described in Application 9.4 for the problem of optimally loading a hopping airplane. Establish a one-to-one correspondence between feasible passenger routings and feasible flows in the minimum cost flow formulation of the problem.
- 9.14.** In this exercise we consider one generalization of the tanker scheduling problem discussed in Application 6.6. Suppose that we can compute the profit associated with each available shipment (depending on the revenues and the operating cost directly attributable to that shipment). Let the profits associated with the shipments 1, 2, 3, and 4 be 10, 10, 3, and 4, respectively. In addition to the operating cost, we incur a fixed charge of 5 units to bring a ship into service. We want to determine the shipments we should make and the number of ships to use to maximize net profits. (Note that it is not necessary to honor all possible shipping commitments.) Formulate this problem as a minimum cost flow problem.
- 9.15.** Consider the following data, with $n = 4$, for the employment scheduling problem that we discussed in Application 9.6. Formulate this problem as a minimum cost flow problem and solve it by the successive shortest path algorithm.

	1	2	3	4	5
1	—	20	35	50	55
$\{c_{ij}\} = 2$	—	—	15	30	40
3	—	—	—	25	35
4	—	—	—	—	10

i	1	2	3	4
$d(i)$	20	15	30	25

9.16. Figure 9.21(b) shows the optimal solution of the minimum cost flow problem shown in Figure 9.21(a). First, verify that x^* is a feasible flow.

- (a) Draw the residual network $G(x^*)$ and show that it contains no negative cycle.
 (b) Specify a set of node potentials π that together with x^* satisfy the reduced cost optimality conditions. List each arc in the residual network and its reduced cost.

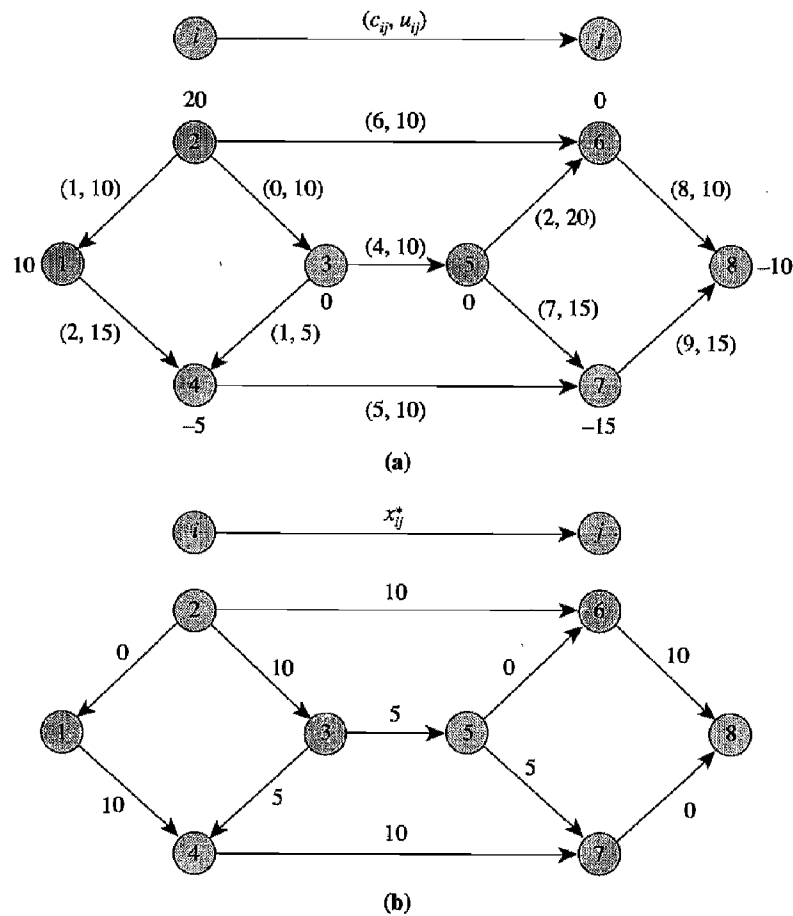


Figure 9.21 Minimum cost flow problem: (a) problem data; (b) optimal solution.

- (c) Verify that the solution x^* satisfies the complementary slackness optimality conditions. To do so, specify a set of optimal node potentials and list the reduced cost of each arc in A .

- 9.17. (a) Figure 9.22(a) gives the data and an optimal solution for a minimum cost flow problem. Assume that all arcs are uncapacitated. Determine optimal node potentials.
- (b) Consider the uncapacitated minimum cost flow problem shown in Figure 9.22(b). For this problem the vector $\pi = (0, -6, -9, -12, -5, -8, -15)$ is an optimal set of node potentials. Determine an optimal flow in the network.

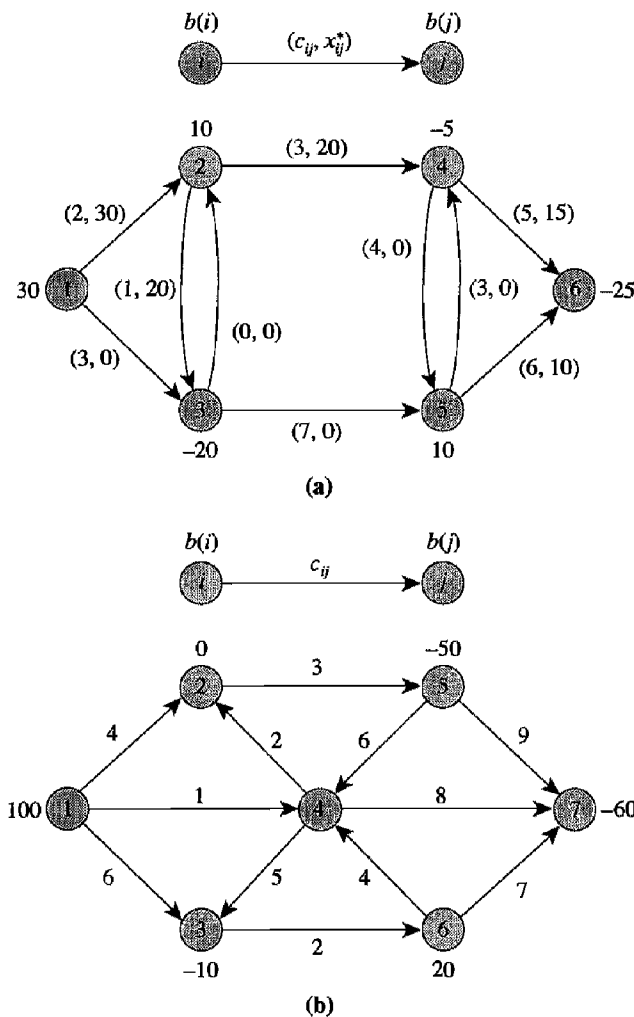


Figure 9.22 Example for Exercise 9.17.

- 9.18. Solve the problem shown in Figure 9.23 by the cycle-canceling algorithm. Use the zero flow as the starting solution.

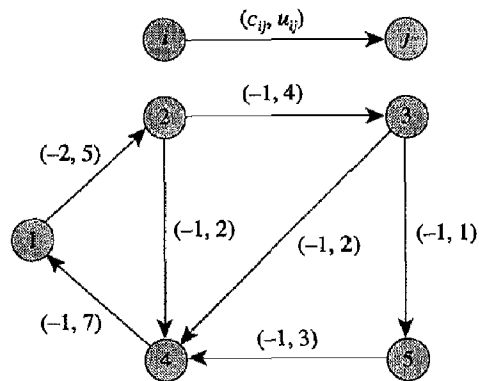


Figure 9.23 Example for Exercise 9.18.

- 9.19.** Show that if we apply the cycle-canceling algorithm to the minimum cost flow problem shown in Figure 9.24, some sequence of augmentations requires 2×10^6 iterations to solve the problem.

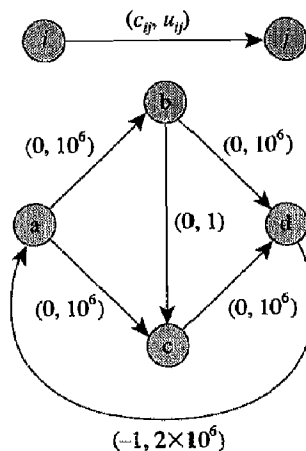


Figure 9.24 Network where cycle canceling algorithm performs 2×10^6 iterations.

- 9.20.** Apply the successive shortest path algorithm to the minimum cost flow problem shown in Figure 9.25. Show that the algorithm performs eight augmentations, each of unit flow, and that the cost of these augmentations (i.e., sum of the arc costs in the path

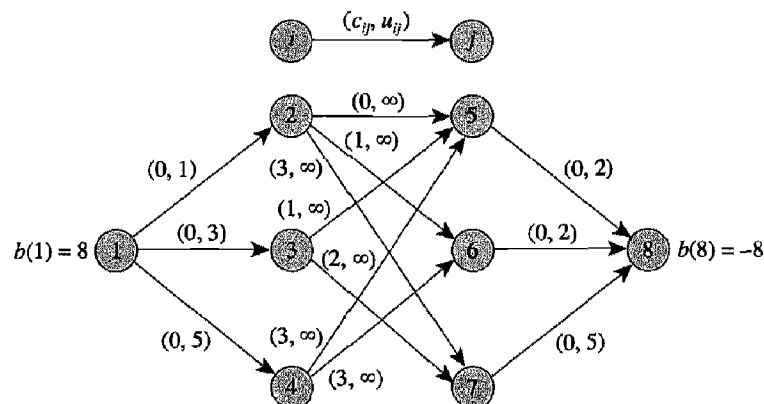


Figure 9.25 Example for Exercise 9.20.

in the residual network) is 0, 1, 2, 3, 3, 4, 5, and 6. How many iterations does the primal–dual algorithm require to solve this problem?

- 9.21. Construct a class of minimum cost flow problems for which the number of iterations performed by the successive shortest path algorithm might grow exponentially in $\log U$. (*Hint*: Consider the example shown in Figure 9.24.)
- 9.22. Figure 9.26 specifies the data and a feasible solution for a minimum cost flow problem. With respect to zero node potentials, list the in-kilter and out-of-kilter arcs. Apply the out-of-kilter algorithm to find an optimal flow in the network.

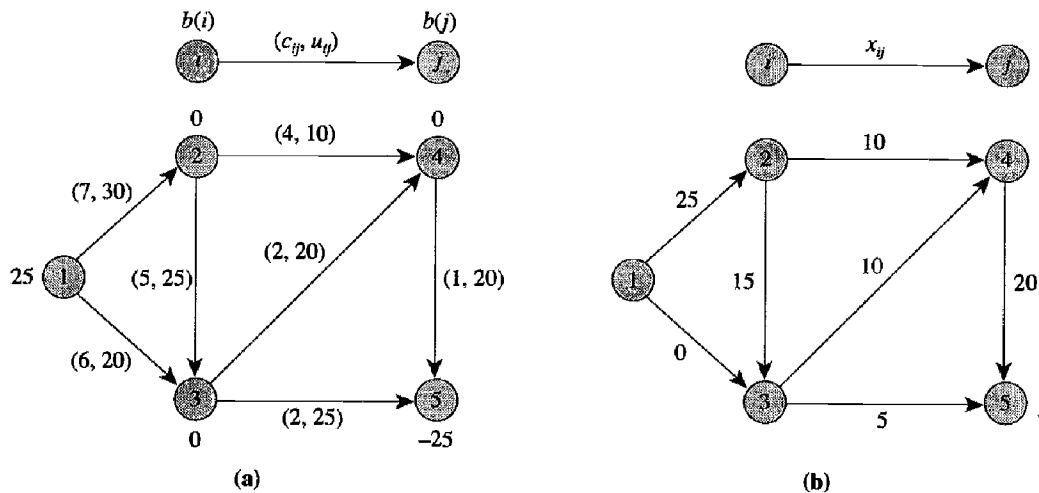


Figure 9.26 Example for Exercises 9.22 and 9.23: (a) problem data; (b) feasible flow.

- 9.23. Consider the minimum cost flow problem shown in Figure 9.26. Starting with zero pseudoflow and zero node potentials, apply the relaxation algorithm to establish an optimal flow.
- 9.24. Figure 9.21(b) specifies an optimal solution for the minimum cost flow problem shown in Figure 9.21(a). Reoptimize the solution with respect to the following changes in the problem data: (1) when c_{23} increases from 0 to 6; (2) when c_{78} decreases from 9 to 2; (3) when $b(2)$ decreases to 15 and $b(8)$ increases to -5 ; and (4) when u_{23} increases to 20. Treat these changes individually.
- 9.25. Assuming that we set one node potential to value zero, show that nC is an upper bound and that $-nC$ is a lower bound on the optimal value of any node potential.
- 9.26. Justify the out-of-kilter algorithm described in Section 9.9 for the case when arcs can violate their flow bounds. Show that in the execution of the algorithm, the kilter number of arcs are nonincreasing and at least one kilter number strictly decreases at every iteration.
- 9.27. Obtain a worst-case bound on the total number of iterations performed by the relaxation algorithm. Compare this bound with the number of iterations performed by the cycle-canceling, successive shortest path, and primal–dual algorithms.
- 9.28. Show that if the pair (x, π) satisfies the complementary slackness optimality conditions (9.8), it also satisfies the reduced cost optimality conditions (9.7).
- 9.29. Prove that if x^* is an optimal flow and π is an optimal set of node potentials, the pair (x^*, π) satisfies the complementary slackness optimality conditions. In your proof, do not use the strong duality theorem. (*Hint*: Suppose that the pair (x, π) satisfies the optimality conditions for some flow x . Show that $c^\pi(x^* - x) = 0$ and use this fact to prove the desired result.)

- 9.30. With respect to an optimal solution x^* of a minimum cost flow problem, suppose that we redefine arc capacities u' as follows:

$$u'_{ij} = \begin{cases} u_{ij} & \text{if } x^*_{ij} = u_{ij} \\ \infty & \text{if } x^*_{ij} < u_{ij}. \end{cases}$$

Show that x^* is also an optimal solution of the minimum cost flow problem with the arc capacities as u' .

- 9.31. With respect to an optimal solution x^* of a minimum cost flow problem, suppose that we redefine arc capacities $u' = x^*$. Show that x^* is also an optimal solution of the minimum cost flow problem with arc capacities u' .
- 9.32. In Section 2.4 we showed how to transform a minimum cost flow problem in an undirected network in which all lower bounds are zero into a minimum cost flow problem in a directed network. Explain why this approach does not work when some lower bounds on arc flows exceed zero.
- 9.33. In the minimum cost flow problem, suppose that one specified arc (p, q) has no lower and upper flow bounds. How would you transform this problem into the standard minimum cost flow problem?
- 9.34. As we have seen in Section 2.4, the uncapacitated transportation problem is equivalent to the minimum cost flow problem in the sense that we can always transform either problem into a version of another problem. If we can solve the uncapacitated transportation problem in $O(g(n, m))$ time, can we also solve the minimum cost flow problem in $O(g(n, m))$ time?
- 9.35. In the *min-cost max-flow problem* defined on a directed network $G = (N, A)$, we wish to send the maximum amount of flow from a node s to a node t at the minimum possible total cost. That is, among all maximum flows, find the one with the smallest cost.
- Show how to formulate any minimum cost flow problem as a min-cost max-flow problem.
 - Show how to convert any min-cost max-flow problem into a circulation problem.
- 9.36. Suppose that in a minimum cost flow problem, some arcs have infinite capacities and some arc costs are negative. (Assume that the lower bounds on all arc flows are zero.)
- Show that the minimum cost flow problem has a finite optimal solution if and only if the uncapacitated arcs do not contain a negative cost-directed cycle.
 - Let B denote the sum of the finite arc capacities and the supplies $b(\cdot)$ of all the supply nodes. Show that the minimum cost flow problem always has an optimal solution in which each arc flow is at most B . Conclude that without any loss of generality, we can assume that in the minimum cost flow problem (with a bounded optimal solution value) every arc is capacitated. (*Hint*: Use the flow decomposition property.)
- 9.37. Suppose that in a minimum cost flow problem, some arcs have infinite capacities and some arc costs are negative. Let B denote the sum of the finite arc capacities and the right-hand-side coefficients $b(i)$ for all the supply nodes. Let z and z' denote the objective function values of the minimum cost flow problem when we set the capacity of each infinite capacity arc to the value B and $B + 1$, respectively. Show that the objective function of the minimum cost flow problem is unbounded if and only if $z' < z$.
- 9.38. In a minimum cost flow network, suppose that in addition to arc capacities, nodes have upper bounds imposed upon the entering flow. Let $w(i)$ be the maximum flow that can enter node $i \in N$. How would you solve this generalization of the minimum cost flow problem?
- 9.39. Let (k, l) and (p, q) denote a minimum cost arc and a maximum cost arc in a network. Is it possible that no minimum cost flow have a positive flow on arc (k, l) ? Is it possible that every minimum cost flow have a positive flow on arc (p, q) ? Justify your answers.
- 9.40. Prove or disprove the following claims.
- Suppose that all supply/demands and arc capacities in a minimum cost flow problem

are all even integers. Then for some optimal flow x^* , each arc flow x_{ij}^* is an even number.

- (b) Suppose that all supply/demands and arc capacities in a minimum cost circulation problem are all even integers. Then for some optimal flow x^* , each arc flow x_{ij}^* is an even number.
- 9.41. Let x^* be an optimal solution of the minimum cost flow problem. Define G° as a subgraph of the residual network $G(x^*)$ consisting of all arcs with zero reduced cost. Show that the minimum cost flow problem has an alternative optimal solution if and only if G° contains a directed cycle.
- 9.42. Suppose that you are given a nonintegral optimal solution to a minimum cost flow problem with integral data. Suggest a method for converting this solution into an integer optimal solution. Your method should maintain optimality of the solution at every step.
- 9.43. Suppose that the pair (x, π) , for some pseudoflow x and some node potentials π , satisfies the reduced cost optimality conditions. Define $G^\circ(x)$ as a subgraph of the residual network $G(x)$ consisting of only those arcs with zero residual capacity. Define the cost of an arc (i, j) in $G^\circ(x)$ as c_{ij} if $(i, j) \in A$, and as $-c_{ij}$ otherwise. Show that every directed path in $G^\circ(x)$ between any pair of nodes is a shortest path in $G(x)$ between the same pair of nodes with respect to the arc costs c_{ij} .
- 9.44. Let x^1 and x^2 be two distinct (alternate) minimum cost flows in a network. Suppose that for some arc (k, l) , $x_{kl}^1 = p$, $x_{kl}^2 = q$, and $p < q$. Show that for every $0 \leq \lambda \leq 1$, the minimum cost flow problem has an optimal solution x (possibly, noninteger) with $x_{kl} = (1 - \lambda)p + \lambda q$.
- 9.45. Let π^1 and π^2 be two distinct (alternate) optimal node potentials of a minimum cost flow problem. Suppose that for some node k , $\pi^1(k) = p$, $\pi^2(k) = q$, and $p < q$. Show that for every $0 \leq \lambda \leq 1$, the minimum cost flow problem has an optimal set of node potentials π (possibly, noninteger) with $\pi(k) = (1 - \lambda)p + \lambda q$.
- 9.46. (a) In the transportation problem, does adding a constant k to the cost of every outgoing arc from a specified supply node affect the optimality of a given optimal solution? Would adding a constant k to the cost of every incoming arc to a specified demand node affect the optimality of a given optimal solution?
- (b) Would your answers to the questions in part (a) be the same if they were posed for the minimum cost flow problem instead of the transportation problem?
- 9.47. In Section 9.7 we described the following practical improvement of the successive shortest path algorithm: (1) terminate the execution of Dijkstra's algorithm whenever it permanently labels a deficit node l , and (2) modify the node potentials by setting $\pi(i)$ to $\pi(i) - d(i)$ if node i is permanently labeled; and by setting $\pi(i)$ to $\pi(i) - d(l)$ if node i is temporarily labeled. Show that after the algorithm has updated the node potentials in this manner, all the arcs in the residual network have nonnegative reduced costs and all the arcs in the shortest path from node k to node l have zero reduced costs. (*Hint*: Use the result in Exercise 5.9.)
- 9.48. Would multiplying each arc cost in a network by a constant k change the set of optimal solutions of the minimum cost flow problem? Would adding a constant k to each arc cost change the set of optimal solutions?
- 9.49. In Section 9.11 we described a method for performing sensitivity analysis when we increase the capacity of an arc (p, q) by 1 unit. Modify the method to perform the analysis when we decrease the capacity of the arc (p, q) by 1 unit.
- 9.50. In Section 9.11 we described a method for performing sensitivity analysis when we increase the cost of an arc (p, q) by 1 unit. Modify the method to perform the analysis when we decrease the cost of the arc (p, q) by 1 unit.
- 9.51. Suppose that we have to solve a minimum cost flow problem in which the sum of the supplies exceeds the sum of the demands, so we need to retain some of the supply at some nodes. We refer to this problem as the *minimum cost flow problem with surplus*. Specify a linear programming formulation of this problem. Also show how to transform this problem into an (ordinary) minimum cost flow problem.

- 9.52.** This exercise concerns the minimum cost flow problem with surplus, as defined in Exercise 9.51. Suppose that we have an optimal solution of a minimum cost flow problem with surplus and we increase the supply of some node by 1 unit, holding the other data fixed. Show that the optimal objective function value cannot increase, but it might decrease. Show that if we increase the demand of a node by 1 unit, holding the other data fixed, the optimal objective function value cannot decrease, but it might increase.
- 9.53. More-for-less paradox** (Charnes and Klingman [1971]). The more-for-less paradox shows that it is possible to send *more* flow from the supply nodes to the demand nodes of a minimum cost flow problem at *lower* cost even if all arc costs are nonnegative. To establish this more-for-less paradox, consider the minimum cost flow problem shown in Figure 9.27. Assume that all arc capacities are infinite.
- (a) Show that the solution given by $x_{14} = 11$, $x_{16} = 9$, $x_{25} = 2$, $x_{26} = 8$, $x_{35} = 11$, and $x_{37} = 14$, is an optimal flow for this minimum cost flow problem. What is the total cost of flow?
- (b) Suppose that we increase the supply of node 2 by 2 units, increase the demand of node 4 by 2 units, and reoptimize the solution using the method described in Section 9.11. Show that the total cost of flow decreases.

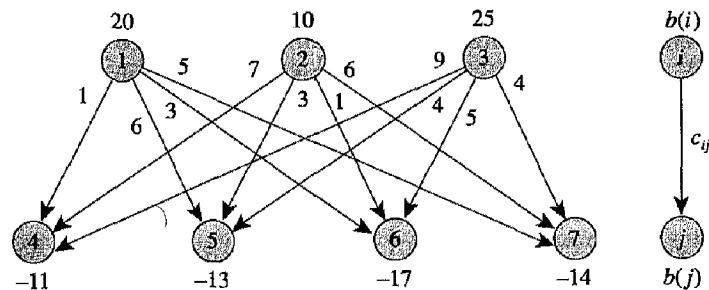


Figure 9.27 More-for-less paradox.

- 9.54. Capacitated minimum spanning tree problem** (Garey and Johnson [1979]). In a complete undirected network with arc lengths c_{ij} and a specially designated node s , called the *central site*, we associate an integer requirement r_i with every node $i \in N - \{s\}$. In the capacitated minimum spanning tree problem, we want to identify a minimum cost spanning tree so that when we send flow on this tree from the central site to the other nodes to satisfy their flow requirements, no tree arc has a flow of more than a given arc capacity R , which is the same for all arcs. Show that when each r_i is 0 or 1, and $R = 1$, we can solve this problem as a minimum cost flow problem. (*Hint*: Model this problem as a minimum cost flow problem with node capacities, as discussed in Exercise 9.38.)
- 9.55. Fractional b -matching problem.** Let $G = (N, A)$ be an undirected graph in which each node $i \in N$ has an associated supply $b(i)$ and each arc $(i, j) \in A$ has an associated cost c_{ij} and capacity u_{ij} . In the b -matching problem, we wish to find a minimum cost subgraph of G with exactly b arcs incident to every node. The *fractional b -matching problem* is a relaxation of the b -matching problem and can be stated as the following linear program:

$$\text{Minimize} \quad \sum_{(i,j) \in A} c_{ij} x_{ij}$$

subject to

$$\sum_{j \in A(i)} x_{ij} = b(i) \quad \text{for all } i \in N,$$

$$0 \leq x_{ij} \leq u_{ij} \quad \text{for all } (i, j) \in A.$$

We assume that $x_{ij} = x_{ji}$ for every arc $(i, j) \in A$. We can define a related minimum cost flow problem as follows. Construct a bipartite network $G' = (N' \cup N'', A')$ with $N' = \{1', 2', \dots, n'\}$, $N'' = \{1'', 2'', \dots, n''\}$, $b(i') = b(i)$, and $b(i'') = -b(i)$. For each arc $(i, j) \in A$, the network G' contains two associated arcs (i', j'') and (j', i'') , each with cost c_{ij} and capacity u_{ij} .

- (a) Show how to transform every solution x of the fractional b -matching problem with cost z into a solution x' of the minimum cost flow problem with cost $2z$. Similarly, show that if x' is a solution of the minimum cost flow problem with cost z' , then $x_{ij} = (x'_{ij} + x'_{ji})/2$ is a feasible solution of the fractional b -matching problem with cost $z'/2$. Use these results to show how to solve the fractional b -matching problem.
 - (b) Show that the fractional b -matching problem always has an optimal solution in which each arc flow x_{ij} is a multiple of $\frac{1}{2}$. Also show that if all the supplies and the capacities are even integers, the fractional b -matching problem always has an integer optimal solution.
- 9.56. Bottleneck transportation problem.** Consider a transportation problem with a traversal time τ_{ij} instead of a cost c_{ij} associated with each arc (i, j) . In the *bottleneck transportation problem* we wish to satisfy the requirements of the demand nodes from the supply nodes in the least time possible [i.e., we wish to find a flow x that minimizes the quantity $\max\{\tau_{ij} : (i, j) \in A \text{ and } x_{ij} > 0\}$].
- (a) Suggest an application of the bottleneck transportation problem.
 - (b) Suppose that we arrange the arc traversal times in the nondecreasing order of their values. Let $\tau_1 < \tau_2 < \dots < \tau_l$ be the distinct values of the arc traversal times (thus $l \leq m$). Let $FS(k, \text{found})$ denote a subroutine that finds whether the transportation problem has a feasible solution using only those the arcs with traversal times less than or equal to τ_k ; assume that the subroutine assigns a value true/false to found. Suggest a method for implementing the subroutine $FS(k, \text{found})$.
 - (c) Using the subroutine $FS(k, \text{found})$, write a pseudocode for solving the bottleneck transportation problem.
- 9.57. Equivalence of minimum cost flow algorithms (Zadeh [1979])**
- (a) Apply the successive shortest path algorithm to the minimum cost flow problem shown in Figure 9.28. Show that it performs four augmentations from node 1 to node 6, each of unit flow.
 - (b) Add the arc $(1, 6)$ with sufficiently large cost and with $u_{16} = 4$ to the example in part (a). Observe that setting $x_{16} = 4$ and $x_{ij} = 0$ for all other arcs gives a feasible flow in the network. With this flow as the initial flow, apply the cycle-canceling algorithm and always augment flow along a negative cycle with minimum cost. Show that this algorithm also performs four unit flow augmentations from node 1 to node 6 along the same paths as in part (a) and in the same order, except that the flow returns to node 1 through the arc $(6, 1)$ in the residual network.

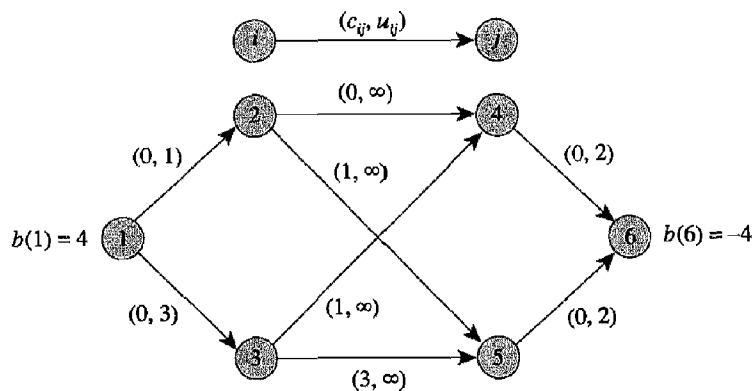


Figure 9.28 Equivalence of minimum cost flow algorithms.

- (c) Using parts (a) and (b) as background, prove the general result that if initialized properly, the successive shortest path algorithm and the cycle-canceling algorithm (with augmentation along a most negative cycle) are equivalent in the sense that they perform the same augmentations and in the same order.
- 9.58.** Modify and initialize the minimum cost flow problem in Figure 9.28 appropriately so that when we apply the out-of-kilter algorithm to this problem, it also performs four augmentation in the same order as the successive shortest path algorithm. Then prove the equivalence of the out-of-kilter algorithm with the successive shortest path algorithm in general.

10

MINIMUM COST FLOWS: POLYNOMIAL ALGORITHMS

Success generally depends upon knowing how long it takes to succeed.
—Montesquieu

Chapter Outline

-
- 10.1 Introduction
 - 10.2 Capacity Scaling Algorithm
 - 10.3 Cost Scaling Algorithm
 - 10.4 Double Scaling Algorithm
 - 10.5 Minimum Mean Cycle-Canceling Algorithm
 - 10.6 Repeated Capacity Scaling Algorithm
 - 10.7 Enhanced Capacity Scaling Algorithm
 - 10.8 Summary
-

10.1 INTRODUCTION

In Chapter 9 we studied several different algorithms for solving minimum cost problems. Although these algorithms guarantee finite convergence whenever the problem data are integral, the computations are not bounded by any polynomial in the underlying problem specification. In the spirit of computational complexity theory, this situation is not completely satisfactory: It does not provide us with any good theoretical assurance that the algorithms will perform well on all problems that we might encounter. The circumstances are quite analogous to our previous development of maximum flow algorithms; we started by first developing straightforward, but not necessarily polynomial, algorithms for solving those problems, and then enhanced these algorithms by changing the algorithmic strategy and/or by using clever data structures and implementations. This situation raises the following natural questions: (1) Can we devise algorithms that are polynomial in the usual problem parameters: number n of nodes, number m of arcs, $\log U$ (the log of the largest supply/demand or arc capacity), and $\log C$ (the log of the largest cost coefficient), and (2) can we develop strongly polynomial-time algorithms (i.e., algorithms whose running time depends upon only on n and m)? A strongly polynomial-time algorithm has one important theoretical advantage: It will solve problems with irrational data.

In this chapter we provide affirmative answers to these questions. To develop polynomial-time algorithms, we use ideas that are similar to those we have used before: namely, scaling of the capacity data and/or of the cost data. We consider

three polynomial-time algorithms: (1) a capacity scaling algorithm that is a scaled version of the successive shortest path algorithm that we discussed in Chapter 9, (2) a cost scaling algorithm that is a generalization of the preflow-push algorithm for the maximum flow problem, and (3) a double scaling algorithm that simultaneously scales both the arc capacities and the costs.

Scaling is a powerful idea that has produced algorithmic improvements to many problems in combinatorial optimization. We might view scaling algorithms as follows. We start with the optimality conditions for the network flow problem we are examining, but instead of enforcing these conditions exactly, we generate an “approximate” solution that is permitted to violate one (or more) of the conditions by an amount Δ . Initially, by choosing Δ quite large, for example as C or U , we will easily be able to find a starting solution that satisfies the relaxed optimality conditions. We then reset the parameter Δ to $\Delta/2$ and reoptimize so that the approximate solution now violates the optimality conditions by an amount of at most $\Delta/2$. We then repeat the procedure, reoptimizing again until the approximate solution violates the conditions by an amount of at most $\Delta/4$, and so on. This solution strategy is quite flexible and leads to different algorithms depending on which of the optimality conditions we relax and how we perform the reoptimizations.

Our discussion of the capacity scaling algorithm for the maximum flow problem in Section 7.3 provides one example. A feasible flow to the maximum flow problem is optimal if the residual network contains no augmenting path. In the capacity scaling algorithm, we relaxed this condition so that after the Δ -scaling phase, the residual network can contain an augmenting path, but only if its capacity were less than Δ . The excess scaling algorithm for the maximum flow problem provides us with another example. In this case the residual network again contains no path from the source node s to the sink node t ; however, at the end of the Δ -scaling phase, we relaxed a feasibility requirement requiring that the flow into every node other than the source and sink equals the flow out of that node. Instead, we permitted the excess at each node to be as large as Δ during the Δ -scaling phase.

In this chapter, by applying a scaling approach to the algorithms that we considered in Chapter 9, we develop polynomial-time versions of these algorithms. We begin by developing a modified version of the successive shortest path algorithm in which we relax two optimality conditions in the Δ -scaling phase: (1) We permit the solution to violate supply/demand constraints by an amount Δ , and (2) we permit the residual network to contain negative cost cycles. The resulting algorithm reduces the number of shortest path computations from nU to $m \log U$.

We next describe a cost-scaling algorithm that uses another concept of approximate optimality; at the end of each ϵ -scaling phase (ϵ plays the role of Δ) we obtain a feasible flow that satisfies the property that the reduced cost of each arc in the residual network is greater than or equal to $-\epsilon$ (instead of zero). To find the optimal solution during the ϵ -scaling phase, this algorithm carries out a sequence of push and relabel operations that are similar to the preflow-push algorithm for maximum flows. The generic cost scaling algorithm runs in $O(n^2 m \log(nC))$ time. We also describe a special “wave implementation” of this algorithm that chooses nodes for the push/relabel operations in a specific order. This specialization requires $O(n^3 \log(nC))$ time.

We then describe a *double scaling algorithm* that combines the features of both cost and capacity scaling. This algorithm works with two nested loops. In the outer loop we scale the costs, and in the inner loop we scale the capacities. Introducing capacity scaling as an inner loop within a cost scaling approach permits us to find augmenting paths very efficiently. This resulting double scaling algorithm solves the minimum cost flow problem in $O(nm \log U \log(nC))$ time.

All of these algorithms require polynomial time; they are not, however, strongly polynomial time because their time bounds depend on $\log U$ and/or $\log C$. Developing strongly polynomial-time algorithms seems to require a somewhat different approach. Although most strongly polynomial-time algorithms use ideas of data scaling, they also use another idea: By invoking the optimality conditions, they are able to show that at intermediate steps of the algorithm, they have already discovered part of the optimal solution (e.g., optimal flow), so that they are able to reduce the problem size. In Sections 10.5, 10.6, and 10.7 we consider three different strongly polynomial-time algorithms whose analysis invokes this “problem reduction argument.”

In Section 10.5 we analyze the minimum mean cycle-canceling algorithm that we described in Section 9.6. Recall that this algorithm augments flow at each step on a cycle with the smallest average cost, averaged over the number of arcs in the cycle, until the residual network contains no negative cost cycle; at this point, the current flow is optimal. As we show in this section, we can view this algorithm as finding a sequence of improved approximately optimal solutions (in the sense that the reduced cost of every arc is greater than or equal to $-\epsilon$, with ϵ decreasing throughout the algorithm). This algorithm has the property that if the magnitude of the reduced cost of any arc is sufficiently large (as a function of ϵ), the flow on that arc remains fixed at its upper or lower bound throughout the remainder of the algorithm and so has this value in the optimal solution. This property permits us to show that the algorithm fixes the flow on an arc and does so sufficiently often so that we obtain an $O(n^2 m^3 \log n)$ time algorithm for the capacitated minimum cost flow problem. One interesting characteristic of this algorithm is that it does not explicitly monitor ϵ or explicitly fix the flow variables. These features of the algorithm are by-products of the analysis.

The strongly polynomial-time algorithm that we consider in Section 10.6 solves the linear programming dual of the minimum cost flow problem. This *repeated capacity scaling algorithm* is a variant of the capacity scaling algorithm that we discuss in Section 10.2. This algorithm uses a scaling parameter Δ as in the capacity scaling algorithm, but shows that periodically the flow on some arc (i, j) becomes sufficiently large (as a function of Δ), at which point we are able to reduce the size of the dual linear program by one, which is equivalent to *contraction* in the primal network. This observation permits us to reduce the size of the problem successively by contracting nodes. The end result is an algorithm requiring $O((m^2 \log n)(m + n \log n))$ time for the minimum cost flow problem.

In Section 10.7 we consider an *enhanced scaling algorithm* that is a hybrid version of the capacity scaling algorithm and the repeated capacity scaling algorithm. By choosing a scaling parameter Δ carefully and by permitting a somewhat broader choice of the augmenting paths at each step, this algorithm is able to fix variables more quickly than the repeated capacity scaling algorithm. As a consequence, it

solves fewer shortest path problems and solves capacitated minimum cost flow problems in $O((m \log n)(m + n \log n))$ time, which is currently the best known polynomial-time bound for solving the capacitated minimum cost flow problem.

10.2 CAPACITY SCALING ALGORITHM

In Chapter 9 we considered the successive shortest path algorithm, one of the fundamental algorithms for solving the minimum cost flow problem. An inherent drawback of this algorithm is that its augmentations might carry relatively small amounts of flow, resulting in a fairly large number of augmentations in the worst case. By incorporating a scaling technique, the capacity algorithm described in this section guarantees that each augmentation carries *sufficiently large* flow and thereby reduces the number of augmentations substantially. This method permits us to improve the worst-case algorithmic performance from $O(nU \cdot S(n, m, nC))$ to $O(m \log U \cdot S(n, m, nC))$. [Recall that U is an upper bound on the largest supply/demand and largest capacity in the network, and $S(n, m, C)$ is the time required to solve a shortest path problem with n nodes, m arcs, and nonnegative costs whose values are no more than C . The reason that the running time involves $S(n, m, nC)$ rather than $S(n, m, C)$ is that the costs in the residual network are reduced costs, and the reduced cost of an arc could be as large as nC .]

The capacity scaling algorithm is a variant of the successive shortest path algorithm. It is related to the successive shortest path algorithm, just as the capacity scaling algorithm for the maximum flow problem (discussed in Section 7.3) is related to the labeling algorithm (discussed in Section 6.5). Recall that the labeling algorithm performs $O(nU)$ augmentations; by sending flows along paths with *sufficiently large* residual capacities, the capacity scaling algorithm reduces the number of augmentations to $O(m \log U)$. In a similar fashion, the capacity scaling algorithm for the minimum cost flow problem ensures that each shortest path augmentation carries a sufficiently large amount of flow; this modification to the algorithm reduces the number of successive shortest path iterations from $O(nU)$ to $O(m \log U)$. This algorithm not only improves on the algorithmic performance of the successive shortest path algorithm, but also illustrates how small changes in an algorithm can produce significant algorithmic improvements (at least in the worst case).

The capacity scaling algorithm applies to the general capacitated minimum cost flow problem. It uses a pseudoflow x and the imbalances $e(i)$ as defined in Section 9.7. The algorithm maintains a pseudoflow satisfying the reduced cost optimality condition and gradually converts this pseudoflow into a flow by identifying shortest paths from nodes with excesses to nodes with deficits and augmenting flows along these paths. It performs a number of scaling phases for different values of a parameter Δ . We refer to a scaling phase with a specific value of Δ as the Δ -scaling phase. Initially, $\Delta = 2^{\lceil \log U \rceil}$. The algorithm ensures that in the Δ -scaling phase each augmentation carries exactly Δ units of flow. When it is not possible to do so because no node has an excess of at least Δ , or no node has a deficit of at least Δ , the algorithm reduces the value of Δ by a factor of 2 and repeats the process. Eventually, $\Delta = 1$ and at the end of this scaling phase, the solution becomes a flow. This flow must be an optimal flow because it satisfies the reduced cost optimality condition.

For a given value of Δ , we define two sets $S(\Delta)$ and $T(\Delta)$ as follows:

$$S(\Delta) = \{i : e(i) \geq \Delta\},$$

$$T(\Delta) = \{i : e(i) \leq -\Delta\}.$$

In the Δ -scaling phase, each augmentation must start at a node in $S(\Delta)$ and end at a node in $T(\Delta)$. Moreover, the augmentation must take place on a path along which every arc has residual capacity of at least Δ . Therefore, we introduce another definition: The Δ -residual network $G(x, \Delta)$ is defined as the subgraph of $G(x)$ consisting of those arcs whose residual capacity is at least Δ . In the Δ -scaling phase, the algorithm augments flow from a node in $S(\Delta)$ to a node in $T(\Delta)$ along a shortest path in $G(x, \Delta)$. The algorithm satisfies the property that every arc in $G(x, \Delta)$ satisfies the reduced cost optimality condition; however, arcs in $G(x)$ but not in $G(x, \Delta)$ might violate the reduced cost optimality condition. Figure 10.1 presents an algorithmic description of the capacity scaling algorithm.

Notice that the capacity scaling algorithm augments exactly Δ units of flow in the Δ -scaling phase, even though it could augment more. For uncapacitated problems, this tactic leads to the useful property that all arc flows are always an integral multiple of Δ . (Why might capacitated networks not satisfy this property?) Several variations of the capacity scaling algorithm discussed in Sections 10.5 and 14.5 adopt the same tactic.

To establish the correctness of the capacity scaling algorithm, observe that the 2Δ -scaling phase ends when $S(2\Delta) = \emptyset$ or $T(2\Delta) = \emptyset$. At that point, either $e(i) < 2\Delta$ for all $i \in N$ or $e(i) > -2\Delta$ for all $i \in N$. These conditions imply that the sum of the excesses (whose magnitude equals the sum of deficits) is bounded by $2n\Delta$.

```

algorithm capacity scaling;
begin
   $x := 0$  and  $\pi := 0$ ;
   $\Delta := 2^{\lceil \log U \rceil}$ ;
  while  $\Delta \geq 1$ 
  begin { $\Delta$ -scaling phase}
    for every arc  $(i, j)$  in the residual network  $G(x)$  do
      if  $r_{ij} \geq \Delta$  and  $c_{ij}^r < 0$  then send  $r_{ij}$  units of flow along arc  $(i, j)$ ,
        update  $x$  and the imbalances  $e(\cdot)$ ;
     $S(\Delta) := \{i \in N : e(i) \geq \Delta\}$ ;
     $T(\Delta) := \{i \in N : e(i) \leq -\Delta\}$ ;
    while  $S(\Delta) \neq \emptyset$  and  $T(\Delta) \neq \emptyset$  do
      begin
        select a node  $k \in S(\Delta)$  and a node  $l \in T(\Delta)$ ;
        determine shortest path distances  $d(\cdot)$  from node  $k$  to all other nodes in the
           $\Delta$ -residual network  $G(x, \Delta)$  with respect to the reduced costs  $c_{ij}^r$ ;
        let  $P$  denote shortest path from node  $k$  to node  $l$  in  $G(x, \Delta)$ ;
        update  $\pi := \pi - d$ ;
        augment  $\Delta$  units of flow along the path  $P$ ;
        update  $x$ ,  $S(\Delta)$ ,  $T(\Delta)$ , and  $G(x, \Delta)$ ;
      end;
       $\Delta := \Delta/2$ ;
    end;
  end;

```

Figure 10.1 Capacity scaling algorithm.

At the beginning of the Δ -scaling phase, the algorithm first checks whether every arc (i, j) in Δ -residual network satisfies the reduced cost optimality condition $c_{ij}^\pi \geq 0$. The arcs introduced in the Δ -residual network at the beginning of the Δ -scaling phase [i.e., those arcs (i, j) for which $\Delta \leq r_{ij} < 2\Delta$] might not satisfy the optimality condition (since, conceivably, $c_{ij}^\pi < 0$). Therefore, the algorithm immediately saturates those arcs (i, j) so that they drop out of the residual network; since the reversal of these arcs (j, i) satisfy the condition $c_{ji}^\pi = -c_{ij}^\pi > 0$, they satisfy the optimality condition. Notice that because $r_{ij} < 2\Delta$, saturating any such arc (i, j) changes the imbalance of its endpoints by at most 2Δ . As a result, after we have saturated all the arcs violating the reduced cost optimality condition, the sum of the excesses is bounded by $2n\Delta + 2m\Delta = 2(n + m)\Delta$.

In the Δ -scaling phase, each augmentation starts at a node $k \in S(\Delta)$, terminates at a node $l \in T(\Delta)$, and carries at least Δ units of flow. Note that Assumption 9.4 implies that the Δ -residual network contains a directed path from node k to node l , so we always succeed in identifying a shortest path from node k to node l . Augmenting flow along a shortest path in $G(x, \Delta)$ preserves the property that every arc satisfies the reduced cost optimality condition (see Section 9.3). When either $S(\Delta)$ or $T(\Delta)$ is empty, the Δ -scaling phase ends. At this point we divide Δ by a factor of 2 and start a new scaling phase. Within $O(\log U)$ scaling phases, $\Delta = 1$, and by the integrality of data, every node imbalance will be zero at the end of this phase. In this phase $G(x, \Delta) \equiv G(x)$ and every arc in the residual network satisfies the reduced cost optimality condition. Consequently, the algorithm will obtain a minimum cost flow at the end of this scaling phase.

As we have seen, the capacity scaling algorithm is easy to state. Similarly, its running time is easy to analyze. We have noted previously that in the Δ -scaling phase the sum of the excesses is bounded by $2(n + m)\Delta$. Since each augmentation in this phase carries at least Δ units of flow from a node in $S(\Delta)$ to a node in $T(\Delta)$, each augmentation reduces the sum of the excesses by at least Δ units. Therefore, a scaling phase can perform at most $2(n + m)$ augmentations. Since we need to solve a shortest path problem to identify each augmenting path, we have established the following result.

Theorem 10.1. *The capacity scaling algorithm solves the minimum cost flow problem in $O(m \log U S(n, m, nC))$ time.* ♦

10.3 COST SCALING ALGORITHM

In this section we describe a cost scaling algorithm for the minimum cost flow problem. This algorithm can be viewed as a generalization of the preflow-push algorithm for the maximum flow problem; in fact, the algorithm reveals an interesting relationship between the maximum flow and minimum cost flow problems. This algorithm relies on the concept of approximate optimality.

Approximate Optimality

A flow x or a pseudoflow x is said to be ϵ -optimal for some $\epsilon > 0$ if for some node potentials π , the pair (x, π) satisfies the following ϵ -optimality conditions:

$$\text{If } c_{ij}^\pi > \epsilon, \text{ then } x_{ij} = 0. \quad (10.1a)$$

$$\text{If } -\epsilon \leq c_{ij}^\pi \leq \epsilon, \text{ then } 0 \leq x_{ij} \leq u_{ij}. \quad (10.1b)$$

$$\text{If } c_{ij}^\pi < -\epsilon, \text{ then } x_{ij} = u_{ij}. \quad (10.1c)$$

These conditions are relaxations of the (exact) complementary slackness optimality conditions (9.8) that we discussed in Section 9.3; note that these conditions reduce to the complementary slackness optimality conditions when $\epsilon = 0$. The *exact* optimality conditions (9.8) imply that any combination of (x_{ij}, c_{ij}^π) lying on the thick lines shown in Figure 10.2(a) is optimal. The ϵ -optimality conditions (10.1) imply that any combination of (x_{ij}, c_{ij}^π) lying on the thick lines or in the hatched region in Figure 10.2(b) is ϵ -optimal.

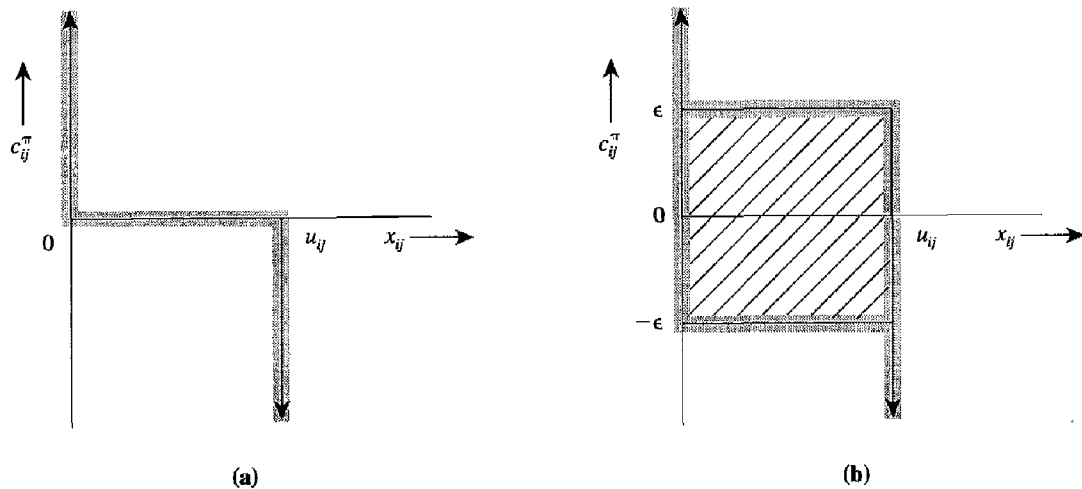


Figure 10.2 Illustrating the optimality condition for arc (i, j) : (a) exact optimality condition for arc (i, j) ; (b) ϵ -optimality condition for arc (i, j) .

The ϵ -optimality conditions assume the following simpler form when stated in terms of the residual network $G(x)$: A flow x or a pseudoflow x is said to be ϵ -optimal for some $\epsilon > 0$ if x , together with some node potential vector π , satisfies the following ϵ -optimality conditions (we leave the proof as an exercise for the reader):

$$c_{ij}^\pi \geq -\epsilon \quad \text{for every arc } (i, j) \text{ in the residual network } G(x). \quad (10.2)$$

Lemma 10.2. *For a minimum cost flow problem with integer costs, any feasible flow is ϵ -optimal whenever $\epsilon \geq C$. Moreover, if $\epsilon < 1/n$, then any ϵ -optimal feasible flow is an optimal flow.*

Proof. Let x be any feasible flow and let $\pi = 0$. Then $c_{ij}^\pi = c_{ij} \geq -C$ for every arc (i, j) in the residual network $G(x)$. Therefore, x is ϵ -optimal for $\epsilon = C$.

Now consider an ϵ -optimal flow x with $\epsilon < 1/n$. Suppose that x is ϵ -optimal with respect to the node potentials π and that W is a directed cycle in $G(x)$. The condition (10.2) implies that $\sum_{(i,j) \in W} c_{ij}^\pi \geq -\epsilon n > -1$, because $\epsilon < 1/n$. The integrality of the costs implies that $\sum_{(i,j) \in W} c_{ij}^\pi$ is nonnegative. But notice that $\sum_{(i,j) \in W} c_{ij}^\pi = \sum_{(i,j) \in W} (c_{ij} - \pi(i) + \pi(j)) = \sum_{(i,j) \in W} c_{ij}$. Therefore, W cannot be a negative cost

cycle. Since $G(x)$ cannot contain any negative cycle, x must be optimal (from Theorem 9.1). ♦

Algorithm

The cost scaling algorithm treats ϵ as a parameter and iteratively obtains ϵ -optimal flows for successively smaller values of ϵ . Initially, $\epsilon = C$ and any feasible flow is ϵ -optimal. The algorithm then performs cost scaling phases by repeatedly applying an *improve-approximation* procedure that transforms an ϵ -optimal flow into a $\frac{1}{2}\epsilon$ -optimal flow. After $1 + \lceil \log(nC) \rceil$ cost scaling phases, $\epsilon < 1/n$ and the algorithm terminates with an optimal flow. Figure 10.3 provides a more formal statement of the cost scaling algorithm.

```

algorithm cost scaling;
begin
   $\pi := 0$  and  $\epsilon := C$ ;
  let  $x$  be any feasible flow;
  while  $\epsilon \geq 1/n$  do
    begin
      improve-approximation( $\epsilon, x, \pi$ );
       $\epsilon := \epsilon/2$ ;
    end;
   $x$  is an optimal flow for the minimum cost flow problem;
end;

```

Figure 10.3 Cost scaling algorithm.

The *improve-approximation* procedure transforms an ϵ -optimal flow into a $\frac{1}{2}\epsilon$ -optimal flow. It does so by (1) converting the ϵ -optimal flow into a $\frac{1}{2}\epsilon$ -optimal pseudoflow, and (2) then gradually converting the pseudoflow into a flow while always maintaining $\frac{1}{2}\epsilon$ -optimality of the solution. We refer to a node i with $e(i) > 0$ as *active* and say that an arc (i, j) in the residual network is *admissible* if $-\frac{1}{2}\epsilon \leq c_{ij} < 0$. The basic operation in the procedure is to select an active node i and perform pushes on admissible arcs (i, j) emanating from node i . When the network contains no admissible arc, the algorithm updates the node potential $\pi(i)$. Figure 10.4 summarizes the essential steps of the generic version of the *improve-approximation* procedure.

Recall that r_{ij} denotes the residual capacity of an arc (i, j) in $G(x)$. As in our earlier discussion of preflow-push algorithms for the maximum flow problem, if $\delta = r_{ij}$, we refer to the push as *saturating*; otherwise, it is *nonsaturating*. We also refer to the updating of the potential of a node as a *relabel* operation. The purpose of a relabel operation at node i is to create new admissible arcs emanating from this node.

We illustrate the basic operations of the *improve-approximation* procedure on a small numerical example. Consider the residual network shown in Figure 10.5(a). Let $\epsilon = 8$. The current pseudoflow is 4-optimal. Node 2 is the only active node in the network, so the algorithm selects it for push/relabel. Suppose that arc $(2, 4)$ is the first admissible arc found. The algorithm pushes $\min\{e(2), r_{24}\} = \min\{30, 5\} = 5$ units of flow on arc $(2, 4)$; this push saturates the arc. Next the algorithm identifies arc $(2, 3)$ as admissible and pushes $\min\{e(2), r_{23}\} = \min\{25, 30\} = 25$ units on this arc. This push is nonsaturating; after the algorithm has performed this push, node

```

procedure improve-approximation( $\epsilon, x, \pi$ );
begin
  for every arc  $(i, j) \in A$  do
    if  $c_{ij}^r > 0$  then  $x_{ij} := 0$ 
    else if  $c_{ij}^r < 0$  then  $x_{ij} := u_{ij}$ ;
  compute node imbalances;
  while the network contains an active node do
    begin
      select an active node  $i$ ;
      push/relabel( $i$ );
    end;
  end;

```

(a)

```

procedure push/relabel( $i$ );
begin
  if  $G(x)$  contains an admissible arc  $(i, j)$  then
    push  $\delta := \min\{e(i), r_{ij}\}$  units of flow from node  $i$  to node  $j$ ;
  else  $\pi(i) := \pi(i) + \epsilon/2$ ;
end;

```

(b)

Figure 10.4 Procedures of the cost scaling algorithm.

2 is inactive and node 3 is active. Figure 10.5(b) shows the residual network at this point.

In the next iteration, the algorithm selects node 3 for push/relabel. Since no admissible arc emanates from this node, we perform a relabel operation and increase the node's potential by $\epsilon/2 = 4$ units. This potential change decreases the reduced costs of the outgoing arcs, namely, (3, 4) and (3, 2), by 4 units and increases the reduced costs of the incoming arcs, namely (1, 3) and (2, 3), by 4 units [see Figure 10.5(c)]. The relabel operation creates an admissible arc, namely arc (3, 4), and we next perform a push of 15 units on this arc [see Figure 10.5(d)]. Since the current solution is a flow, the improve-approximation procedure terminates.

To identify admissible arcs emanating from node i , we use the same data structure used in the maximum flow algorithms described in Section 7.4. For each node i , we maintain a *current-arc* (i, j) which is the current candidate to test for admissibility. Initially, the current-arc of node i is the first arc in its arc list $A(i)$. To determine an admissible arc emanating from node i , the algorithm checks whether the node's current-arc is admissible, and if not, chooses the next arc in the arc list as the current-arc. Thus the algorithm passes through the arc list starting with the current-arc until it finds an admissible arc. If the algorithm reaches the end of the arc list without finding an admissible arc, it declares that the node has no admissible arc. At this point it relabels node i and again sets its current-arc to the first arc in $A(i)$.

We might comment on two practical improvements of the improve-approximation procedure. The algorithm, as stated, starts with $\epsilon = C$ and reduces ϵ by a factor of 2 in every scaling phase until $\epsilon = 1/n$. As a consequence, ϵ could become

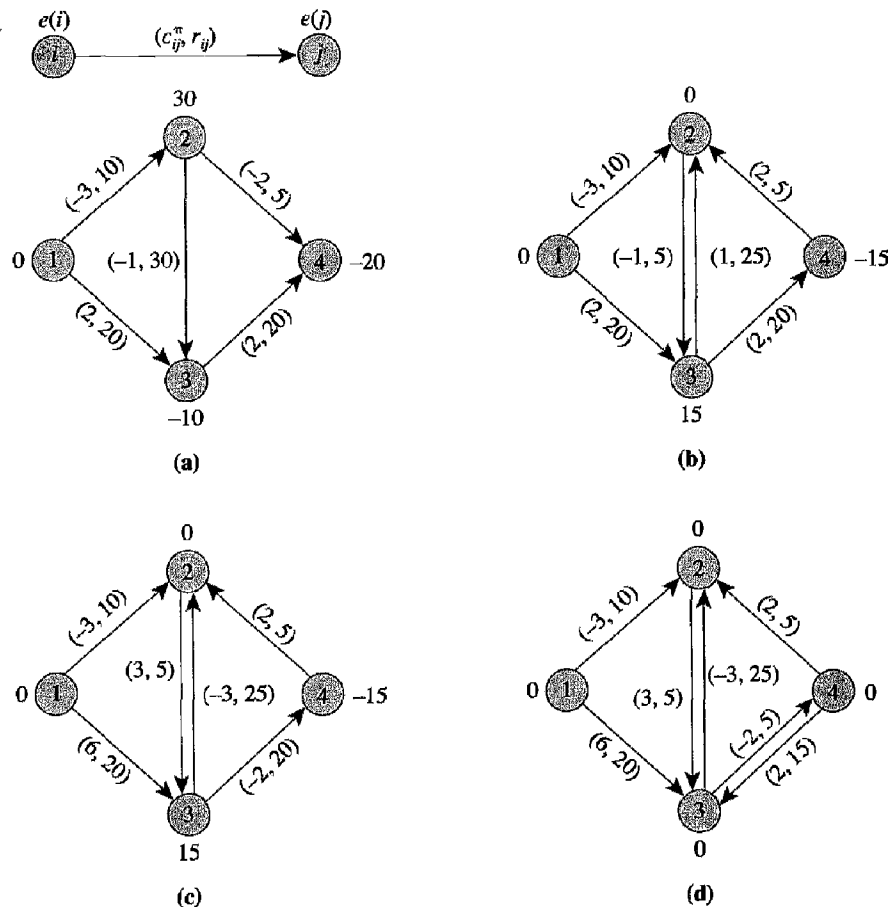


Figure 10.5 Illustration of push/relabel steps.

nonintegral during the execution of the algorithm. By slightly modifying the algorithm, however, we can ensure that ϵ remains integral. We do so by multiplying all the arc costs by n , by setting the initial value of ϵ equal to $2^{\lceil \log(nC) \rceil}$, and by terminating the algorithm when $\epsilon < 1$. It is possible to show (see Exercise 10.7) that the modified algorithm would yield an optimal flow for the minimum cost flow problem in the same computational time. Furthermore, as stated, the algorithm increases a node potential by $\epsilon/2$ during a relabel operation. As described in Exercise 10.8, we can often increase the node potential by an amount larger than $\epsilon/2$.

Analysis of the Algorithm

We show that the cost scaling algorithm correctly solves the minimum cost flow problem. In the proof, we rely on the fact that the improve-approximation procedure converts an ϵ -optimal flow into an $\epsilon/2$ -optimal flow. We establish this result in the following lemma.

Lemma 10.3. *The improve-approximation procedure always maintains $\frac{1}{2} \epsilon$ -optimality of the pseudoflow, and at termination yields a $\frac{1}{2} \epsilon$ -optimal flow.*

Proof. We use induction on the number of pushes and relabels. At the beginning of the procedure, the algorithm sets the flows on arcs with negative reduced costs to their capacities, sets the flow on arcs with positive reduced costs to zero, and leaves the flow on arcs with zero reduced costs unchanged. The resulting pseudoflow satisfies (10.1) for $\epsilon = 0$ and thus is 0-optimal. Since a 0-optimal pseudoflow is ϵ -optimal for every ϵ , the resulting flow is also $\frac{1}{2}\epsilon$ -optimal.

We next study the effect of a push on the $\frac{1}{2}\epsilon$ -optimality of the solution. Pushing flow on arc (i, j) might add its reversal (j, i) to the residual network. But since $-\epsilon/2 \leq c_{ij}^\pi < 0$ (by the criteria of admissibility), $c_{ji}^\pi = -c_{ij}^\pi > 0$, and so this arc satisfies the $\frac{1}{2}\epsilon$ -optimality condition (10.2).

What is the effect of a relabel operation? The algorithm relabels a node i when $c_{ij}^\pi \geq 0$ for every arc (i, j) emanating from node i in the residual network. Increasing the potential of node i by $\epsilon/2$ units decreases the reduced cost of all arcs emanating from node i by $\epsilon/2$ units. But since $c_{ij}^\pi \geq 0$ before the increase in π , $c_{ij}^\pi \geq -\epsilon/2$ after the increase, and the arc satisfies the $\frac{1}{2}\epsilon$ -optimality condition. Furthermore, increasing the potential of node i by $\epsilon/2$ units increases the reduced costs of the incoming arcs at node i but maintains the $\frac{1}{2}\epsilon$ -optimality condition for these arcs. These results establish the lemma. ♦

We next analyze the complexity of the improve-approximation procedure. We show that the number of relabel operations is $O(n^2)$, the number of saturating pushes is $O(nm)$, and the number of nonsaturating pushes for the generic version is $O(n^2m)$. These time bounds are comparable to those of the preflow-push algorithms for the maximum flow problem and the proof techniques are also similar. We first prove the most significant result, which bounds the number of relabel operations.

Lemma 10.4. *No node potential increases more than $3n$ times during an execution of the improve-approximation procedure.*

Proof. Let x be the current $\frac{1}{2}\epsilon$ -optimal pseudoflow and x' be the ϵ -optimal flow at the end of the previous cost scaling phase. Let π and π' be the node potentials corresponding to the pseudoflow x and the flow x' . It is possible to show (see Exercise 10.9) that for every node v with an excess there exists a node w with a deficit and a sequence of nodes $v = v_0, v_1, v_2, \dots, v_l = w$ that satisfies the property that the path $P = v_0 - v_1 - v_2 - \dots - v_l$ is a directed path in $G(x)$ and its reversal $\bar{P} = v_l - v_{l-1} - \dots - v_1 - v_0$ is a directed path in $G(x')$. Applying the $\frac{1}{2}\epsilon$ -optimality condition to the arcs on the path P in $G(x)$, we see that

$$\sum_{(i,j) \in P} c_{ij}^\pi \geq -l(\epsilon/2).$$

Substituting $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$ in this expression gives

$$\sum_{(i,j) \in P} c_{ij} - \pi(v) + \pi(w) \geq -l(\epsilon/2).$$

Alternatively,

$$\pi(v) \leq \pi(w) + l(\epsilon/2) + \sum_{(i,j) \in P} c_{ij}. \quad (10.3)$$

Applying the ϵ -optimality conditions to the arcs on the path \bar{P} in $G(x')$, we obtain $\sum_{(j,i) \in \bar{P}} c_{ji}' \geq -l\epsilon$. Substituting $c_{ji}' = c_{ji} - \pi'(j) + \pi'(i)$ in this expression gives

$$\sum_{(j,i) \in \bar{P}} c_{ji} - \pi'(w) + \pi'(v) \geq -l\epsilon. \quad (10.4)$$

Notice that $\sum_{(j,i) \in \bar{P}} c_{ji} = -\sum_{(i,j) \in P} c_{ij}$ since \bar{P} is a reversal of P . In view of this fact, we can restate (10.4) as

$$\sum_{(i,j) \in P} c_{ij} \leq l\epsilon - \pi'(w) + \pi'(v). \quad (10.5)$$

Substituting (10.5) in (10.3), we see that

$$(\pi(v) - \pi'(v)) \leq (\pi(w) - \pi'(w)) + 3l\epsilon/2. \quad (10.6)$$

Now we use the facts that (1) $\pi(w) = \pi'(w)$ (the potential of a node with negative imbalance does not change because the algorithm never selects it for push/relabel), (2) $l \leq n$, and (3) each increase in the potential increases $\pi(v)$ by at least $\epsilon/2$ units. These facts and expression (10.6) establish the lemma. ♦

Lemma 10.5. *The improve-approximation procedure performs $O(nm)$ saturating pushes.*

Proof. We show that between two consecutive saturations of an arc (i, j) , the procedure must increase both the potentials $\pi(i)$ and $\pi(j)$ at least once. Consider a saturating push on arc (i, j) . Since arc (i, j) is admissible at the time of the push, $c_{ij}^\pi < 0$. Before the algorithm can saturate this arc again, it must send some flow back from node j to node i . At that time $c_{ji}^\pi < 0$ or $c_{ij}^\pi > 0$. These conditions are possible only if the algorithm has relabeled node j . In the subsequent saturation of arc (i, j) , $c_{ij}^\pi < 0$, which is possible only if the algorithm has relabeled node i . But by the previous lemma the improve-approximation procedure can relabel any node $O(n)$ times, so it can saturate any arc $O(n)$ times. Consequently, the number of saturating pushes is $O(nm)$. ♦

To bound the number of nonsaturating pushes, we need one more result. We define the *admissible network* of a given residual network as the network consisting solely of admissible arcs. For example, Figure 10.6(b) specifies the admissible network for the residual network given in Figure 10.6(a).

Lemma 10.6. *The admissible network is acyclic throughout the improve-approximation procedure.*

Proof. We show that the algorithm satisfies this property at every step. The result is true at the beginning of the improve-approximation procedure because the initial pseudoflow is 0-optimal and the residual network contains no admissible arc. We show that the result remains valid throughout the procedure. We always push flow on arc (i, j) with $c_{ij}^\pi < 0$; therefore, if the algorithm adds the reversal (j, i) of this arc to the residual network, then $c_{ji}^\pi > 0$ and so the reversal arc is nonadmissible. Thus pushes do not create new admissible arcs and the admissible network remains acyclic. The relabel operation at node i decreases the reduced costs of all outgoing

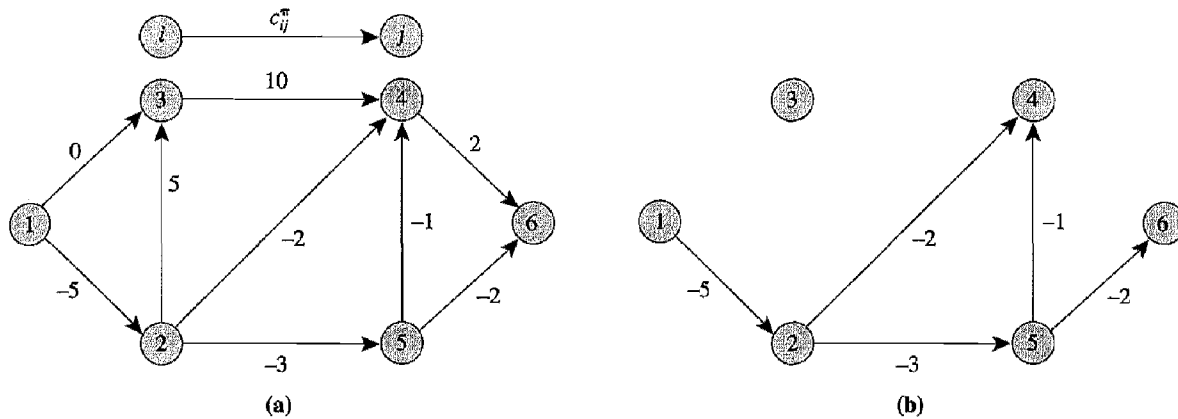


Figure 10.6 Illustration of an admissible network: (a) residual network; (b) admissible network.

arcs at node i by $\epsilon/2$ units and might create new admissible arcs. This relabel operation increases the reduced costs of all incoming arcs at node i by $\epsilon/2$ units, so all such arcs become inadmissible. Consequently, the relabel operation cannot create any directed cycle passing through node i . Thus neither of the two operations, pushes and relabels, of the algorithm can create a directed cycle, which establishes the lemma. \blacklozenge

Lemma 10.7. *The improve-approximation procedure performs $O(n^2m)$ non-saturating pushes.*

Proof. We use a potential function argument to prove the lemma. Let $g(i)$ be the number of nodes that are reachable from node i in the admissible network and let $\Phi = \sum_{i \text{ is active}} g(i)$ be the potential function. We assume that every node is reachable from itself. For example, in the admissible network shown in Figure 10.7, nodes 1 and 4 are the only active nodes. In this network, nodes 1, 2, 3, 4, and 5 are reachable from node 1, and nodes 4 and 5 are reachable from node 4. Therefore, $g(1) = 5$, $g(4) = 2$, and $\Phi = 7$.

At the beginning of the procedure, $\Phi \leq n$ since the admissible network contains

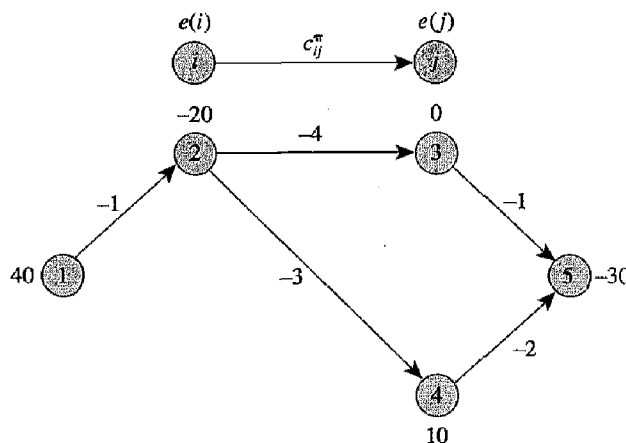


Figure 10.7 Admissible network for $\epsilon = 8$.

no arc and each $g(i) = 1$. After a saturating push on arc (i, j) , node j might change its state from inactive to active, which would increase Φ by $g(j) \leq n$. Therefore, Lemma 10.5 implies that the total increase due to saturating pushes is $O(n^2m)$. A relabel operation of node i might create new admissible arcs (i, j) and will increase $g(i)$ by at most n units. But this relabel operation does not increase $g(k)$ for any other node k because it makes all incoming arcs at node k inadmissible (see the proof of Lemma 10.6). Thus the total increase due to all relabel operations is $O(n^3)$.

Finally, consider the effect on Φ of a nonsaturating push on arc (i, j) . As a result of the push, node i becomes inactive and node j might change its status from inactive to active. Thus the push decreases Φ by $g(i)$ units and might increase it by another $g(j)$ units. Now notice that $g(i) \geq g(j) + 1$ because every node that is reachable from node j is also reachable from node i but node i is not reachable from node j (because the admissible network is acyclic). Therefore, a nonsaturating push decreases Φ by at least 1 unit. Consequently, the total number of nonsaturating pushes is bounded by the initial value of Φ plus the total increase in Φ throughout the algorithm, which is $O(n) + O(n^2m) + O(n^3) = O(n^2m)$. This result establishes the lemma. ♦

Let us summarize our discussion. The improve-approximation procedure requires $O(n^2m)$ time to perform nonsaturating pushes and $O(nm)$ time to perform saturating pushes. The amount of time needed to identify admissible arcs is $O(\sum_{i \in N} |A(i)|n) = O(nm)$, since between two consecutive potential increases of a node i , the algorithm will examine $|A(i)|$ arcs for testing admissibility. The algorithm could store all the active nodes in a list. Doing so would permit it to identify an active node in $O(1)$ time, so this operation would not be a bottleneck step. Consequently, the improve-approximation procedure runs in $O(n^2m)$ time. Since the cost scaling algorithm calls this procedure $1 + \lceil \log(nC) \rceil$ times, we obtain the following result.

Theorem 10.8. *The generic cost scaling algorithm runs in $O(n^2m \log(nC))$ time.* ♦

The cost scaling algorithm illustrates an important connection between the maximum flow and the minimum cost flow problems. Solving an improve-approximation problem is very similar to solving a maximum flow problem by the preflow-push method. Just as in the preflow-push algorithm, the bottleneck operation in the procedure is the number of nonsaturating pushes. In Chapter 7 we have seen how to reduce the number of nonsaturating pushes for the preflow-push algorithm by examining active nodes in some specific order. Similar ideas permit us to streamline the improve-approximation procedure as well. We describe one such improvement, called the *wave implementation*, that reduces the number of nonsaturating pushes from $O(n^2m)$ to $O(n^3)$.

Wave Implementation

Before we describe the wave implementation, we introduce the concept of *node examination*. In an iteration of the improve-approximation procedure, the algorithm selects a node, say node i , and either performs a saturating push or a nonsaturating

push from this node, or relabels the node. If the algorithm performs a saturating push, node i might still be active, but the algorithm might select another node in the next iteration. We shall henceforth assume that whenever the algorithm selects a node, it keeps pushing flow from that node until either its excess becomes zero or the node becomes relabeled. If we adopt this node selection strategy, the algorithm will perform several saturating pushes from a particular node followed either by a nonsaturating push or a relabel operation; we refer to this sequence of operations as a *node examination*.

The wave implementation is a special implementation of the improve-approximation procedure that selects active nodes for push/relabel steps in a specific order. The algorithm uses the fact that the admissible network is acyclic. In Section 3.4 we showed that it is always possible to order nodes of an acyclic network so that for every arc (i, j) in the network, node i occurs prior to node j . Such an ordering of nodes is called a *topological ordering*. For example, for the admissible network shown in Figure 10.6, one possible topological ordering of nodes is 1–2–5–4–3–6. In Section 3.4 we showed how to arrange the nodes of a network in a topological order in $O(m)$ time. For a given topological order, we define the *rank* of a node as n minus its number in the topological sequence. For example, in the preceding example, $\text{rank}(1) = 6$, $\text{rank}(6) = 1$ and $\text{rank}(5) = 4$.

Observe that each push carries flow from a node with higher rank to a node with lower rank. Also observe that pushes do not change the topological ordering of nodes since they do not create new admissible arcs. The relabel operations, however, might create new admissible arcs and consequently, might affect the topological ordering of nodes.

The wave implementation sequentially examines nodes in the topological order and if the node being examined is active, it performs push/relabel steps at the node until either the node becomes inactive or it becomes relabeled. When examined in this order, the active nodes push their excesses to nodes with lower rank, which in turn push their excesses to nodes with even lower rank, and so on. A relabel operation changes the topological order; so after each relabel operation the algorithm modifies the topological order and again starts to examine nodes according to the topological order. If within n consecutive node examinations, the algorithm performs no relabel operation, then at this point all the active nodes have discharged their excesses and the algorithm has obtained a flow. Since the algorithm performs $O(n^2)$ relabel operations, we immediately obtain a bound of $O(n^3)$ on the number of node examinations. Each node examination entails at most one nonsaturating push. Consequently, the wave algorithm performs $O(n^3)$ nonsaturating pushes per execution of improve-approximation.

To illustrate the wave implementation, we consider the pseudoflow shown in Figure 10.8. One topological order of nodes is 2–3–4–1–5–6. The algorithm first examines node 2 and pushes 20 units of flow on arc $(2, 1)$. Then it examines node 3 and pushes 5 units of flow on arc $(3, 1)$ and 10 units of flow on arc $(3, 4)$. The push creates an excess of 10 units at node 4. Next the algorithm examines node 4 and sends 5 units on the arc $(4, 6)$. Since node 4 has an excess of 5 units but has no outgoing admissible arc, we need to relabel node 4 and reexamine all nodes in the topological order starting with the first node in the order.

To complete the description of the algorithm, we need to describe a procedure

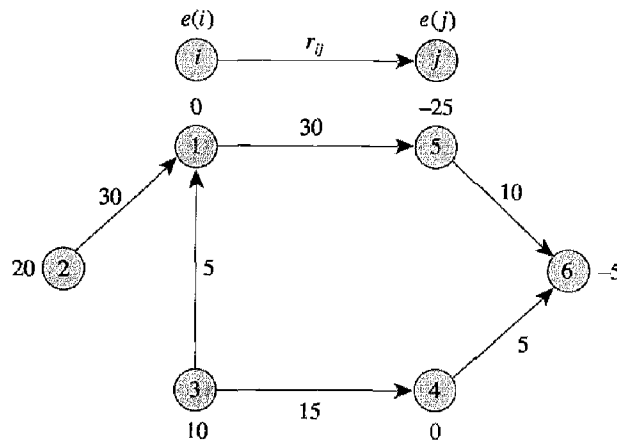


Figure 10.8 Example to illustrate the wave implementation.

for obtaining a topological order of nodes after each relabel operation. We can use an $O(m)$ algorithm to determine an initial topological ordering of the nodes (see Section 3.4). Suppose that while examining node i , the algorithm relabels this node. At this point, the network contains no incoming admissible arc at node i . We claim that if we move node i from its present position to the first position in the previous topological order leaving all other nodes intact, we obtain a topological order of the new admissible network. For example, for the admissible network given in Figure 10.8, one topological order of the nodes is 2–3–4–1–5–6. If we examine nodes in this order, the algorithm relabels node 4. After the algorithm has performed this relabel operation, the modified topological order of nodes is 4–2–3–1–5–6. This method works because (1) after the relabeling, node i has no incoming admissible arc, so assigning it to the first place in the topological order is justified; (2) the relabeling of node i might create some new outgoing admissible arcs (i, j) but since node i is first in the topological order, any such arc satisfies the conditions of a topological ordering; and (3) the rest of the admissible network does not change, so the previous order remains valid. Therefore, the algorithm maintains an ordered set of nodes (possibly as a doubly linked list) and examines nodes in this order. Whenever it relabels a node i , the algorithm moves this node to the first place in the order and again examines nodes in order starting from node i .

We have established the following result.

Theorem 10.9. *The wave implementation of the cost scaling algorithm solves the minimum cost flow problem in $O(n^3 \log(nC))$ time.* ♦

By examining the active nodes carefully and thereby reducing the number of nonsaturating pushes, the wave implementation improves the running time of the generic implementation of the improve-approximation procedure from $O(n^2m)$ to $O(n^3)$. A complementary approach for improving the running time is to use cleverer data structure to reduce the time per nonsaturating push. Using the *dynamic tree* data structures described in Section 8.5, we can improve the running time of the generic implementation to $O(nm \log n)$ and of the wave implementation to $O(nm \log(n^2/m))$. The references cited at the end of the chapter contain the details of these implementations.

10.4 DOUBLE SCALING ALGORITHM

As we have seen in the preceding two sections, by scaling either the arc capacities or the cost coefficients of a minimum cost flow problem, we can devise algorithms with improved worst-case performance. This development raises a natural question: Can we combine ideas from these algorithms to obtain even further improvements that are not obtained by either technique alone? In this section we provide an affirmative answer to this question. The double scaling algorithm we describe solves the capacitated minimum cost flow problem in $O(nm \log U \log(nC))$ time. When implemented using a dynamic tree data structure, this approach produces one of the best polynomial time algorithms for solving the minimum cost flow problem.

In this discussion we assume that the reader is familiar with the capacity scaling algorithm and the cost scaling algorithm that we examined in the preceding two sections. To solve the capacitated minimum cost flow problem, we first transform it into an uncapacitated transportation problem using the transformation described in Section 2.4. We assume that every arc in the minimum cost flow problem is capacitated. Consequently, the transformed network will be a bipartite network $G = (N_1 \cup N_2, A)$ with N_1 and N_2 as the sets of supply and demand nodes. Moreover, $|N_1| = n$ and $|N_2| = m$.

The double scaling algorithm is the same as the cost scaling algorithm described in the preceding section except that it uses a more efficient version of the improve-approximation procedure. The improve-approximation procedure in the preceding section relied on a “pseudoflow-push” method to push flow out of active nodes. A natural alternative would be to try an augmenting path based method. This approach would send flow from a node with excess to a node with deficit over an *admissible path* (i.e., a path in which each arc is admissible). A straightforward implementation of this approach would require $O(nm)$ augmentations since each augmentation would saturate at least one arc and, by Lemma 10.5, the algorithm requires $O(nm)$ arc saturations. Since each augmentation requires $O(n)$ time, this approach does not appear to improve the $O(n^2m)$ bound of the generic improve-approximation procedure.

We can, however, use ideas from the capacity scaling algorithm to reduce the number of augmentations to $O(m \log U)$ by ensuring that each augmentation carries *sufficiently large* flow. The resulting algorithm performs cost scaling in an “outer loop” to obtain ϵ -optimal flows for successively smaller values of ϵ . Within each cost scaling phase, we start with a pseudoflow and perform a number of capacity scaling phases, called Δ -scaling phases, for successively smaller values of Δ . In the Δ -scaling phase, the algorithm identifies admissible paths from a node with an excess of at least Δ to a node with a deficit and augments Δ units of flow over these paths. When all node excesses are less than Δ , we reduce Δ by a factor of 2 and initiate a new Δ -scaling phase. At the end of the 1-scaling phase, we obtain a flow.

The algorithmic description of the double scaling algorithm is same as that of the cost scaling algorithm except that we replace the improve-approximation procedure by the procedure given in Figure 10.9.

The capacity scaling within the *improve-approximation* procedure is somewhat different from the capacity scaling algorithm described in Section 10.2. The new algorithm differs from the one we considered previously in the following respects:

```

procedure improve-approximation( $\epsilon, x, \pi$ );
begin
  set  $x := 0$  and compute node imbalances;
   $\pi(j) := \pi(j) + \epsilon$ , for all  $j \in N_2$ ;
   $\Delta := 2^{\lceil \log U \rceil}$ ;
  while the network contains an active node do
    begin
       $S(\Delta) := \{i \in N_1 \cup N_2 : e(i) \geq \Delta\}$ ;
      while  $S(\Delta) \neq \emptyset$  do
        begin  $\{\Delta$ -scaling phase $\}$ 
          select a node  $k$  from  $S(\Delta)$ ;
          determine an admissible path  $P$  from node  $k$  to some node  $l$  with  $e(l) < 0$ ;
          augment  $\Delta$  units of flow on path  $P$  and update  $x$  and  $S(\Delta)$ ;
        end;
         $\Delta := \Delta/2$ ;
      end;
    end;
end;

```

Figure 10.9 Improve-approximation procedure in the double scaling algorithm.

(1) the augmentation terminates at a node l with $e(l) < 0$ but whose deficit may not be as large as Δ ; (2) each residual capacity is an integral multiple of Δ because each arc flow is an integral multiple of Δ and each arc capacity is ∞ ; and (3) the algorithm does not change flow on some arcs at the beginning of the Δ -scaling phase to ensure that the solution satisfies the optimality conditions. We point out that the algorithm feature (3) is a consequence of feature (2) because each r_{ij} is a multiple of Δ , so $G(x, \Delta) \equiv G(x)$.

The double scaling algorithm improves on the capacity scaling algorithm by identifying an admissible path in only $O(n)$ time, on average, rather than the time $O(S(n, m, nC))$ required to identify an augmenting path in the capacity scaling algorithm. The savings in identifying augmenting paths more than offsets the extra requirement of performing $O(\log(nC))$ cost scaling phases in the double scaling algorithm.

We next describe a method for identifying admissible paths efficiently. The algorithm identifies an admissible path by starting at node k and gradually building up the path. It maintains a *partial admissible path* P , which is initially null, and keeps enlarging it until it includes a node with deficit. We maintain the partial admissible path P using predecessor indices [i.e., if $(u, v) \in P$ then $\text{pred}(v) = u$]. At any point in the algorithm, we perform one of the following two steps, whichever is applicable, from the tip of P (say, node i):

advance(i). If the residual network contains an admissible arc (i, j) , add (i, j) to P and set $\text{pred}(j) := i$. If $e(j) < 0$, stop.

retreat(i). If the residual network does not contain an admissible arc (i, j) , update $\pi(i)$ to $\pi(i) + \epsilon/2$. If $i \neq k$, remove the arc $(\text{pred}(i), i)$ from P so that $\text{pred}(i)$ becomes its new tip.

The retreat step relabels (increases the potential of) node i for the purpose of creating new admissible arcs emanating from this node. However, increasing the potential of node i increases the reduced costs of all the incoming arcs at the node

i by $\epsilon/2$. Consequently, the arc $(\text{pred}(i), i)$ becomes inadmissible, so we delete this arc from P (provided that P is nonempty).

We illustrate the method for identifying admissible paths on the example shown in Figure 10.10. Let $\epsilon = 4$ and $\Delta = 4$. Since node 1 is the only node with an excess of at least 4, we begin to develop the admissible path starting from this node. We perform the step $\text{advance}(1)$ and add the arc $(1, 2)$ to P . Next, we perform the step $\text{advance}(2)$ and add the arc $(2, 4)$ to P . Now node 4 has no admissible arc. So we perform a retreat step. We increase the potential of node 4 by $\epsilon/2 = 2$ units, thus changing the reduced cost of arc $(2, 4)$ to 1; so we eliminate this arc from P . In the next two steps, the algorithm performs the steps $\text{advance}(2)$ and $\text{advance}(5)$, adding arcs $(2, 5)$ and $(5, 6)$ to P . Since the path now contains node 6, which is a node with a deficit, the method terminates. It has found the admissible path 1–2–5–6.

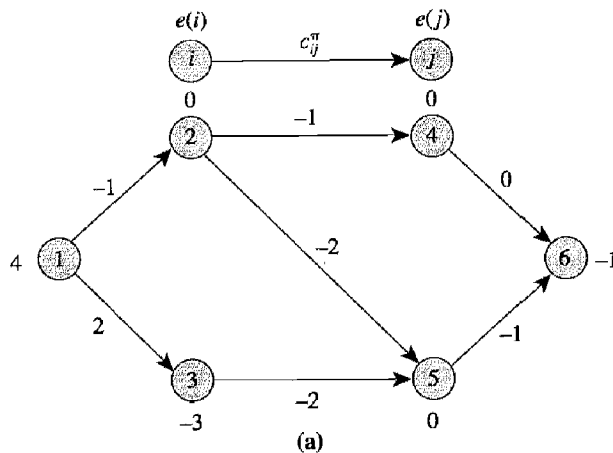


Figure 10.10 Residual network.

It is easy to show that the double scaling algorithm correctly solves the minimum cost flow problem. At the beginning of the improve-approximation procedure, we set $x = 0$ and the corresponding residual network is the same as the original network. The ϵ -optimality of the solution at the end of the previous scaling phase implies that $c_{ij}^\pi \geq -\epsilon$ for all arcs $(i, j) \in A$. Therefore, by adding ϵ to $\pi(j)$ for each $j \in N_2$, we obtain an $\frac{1}{2}\epsilon$ -optimal pseudoflow (in fact, it is a 0-optimal pseudoflow). Like the improve-approximation procedure described in the preceding section, the algorithm always augments flow on admissible arcs and relabels a node when it has no outgoing admissible arc. Consequently, the algorithm preserves $\frac{1}{2}\epsilon$ -optimality of the pseudoflow and at termination yields a $\frac{1}{2}\epsilon$ -optimal flow.

We next consider the complexity of the improve-approximation procedure. Each execution of the procedure performs $(1 + \lfloor \log U \rfloor)$ capacity scaling phases. At the end of the 2Δ -scaling phase, $S(2\Delta) = \emptyset$. Therefore, at the beginning of the Δ -scaling phase, $\Delta \leq e(i) < 2\Delta$ for each node $i \in S(\Delta)$. During the Δ -scaling phase, the algorithm augments Δ units of flow from a node k in $S(\Delta)$ to a node l with $e(l) < 0$. The augmentation reduces the excess of node k to a value less than Δ and ensures that the imbalance at node l is strictly less than Δ . Consequently, each augmentation deletes a node from $S(\Delta)$ and after at most $|N_1| + |N_2| = O(m)$ augmentations, $S(\Delta)$ becomes empty and the algorithm begins a new capacity scaling phase. The algorithm thus performs a total of $O(m \log U)$ augmentations.

We next focus on the time needed to identify admissible paths. We first count the number of advance steps. Each advance step adds an arc to the partial admissible path, and each retreat step deletes an arc from the partial admissible path. Thus we can distinguish between two types of advance steps: (1) those that add arcs to an admissible path on which the algorithm later performs an augmentation, and (2) those that are later canceled by a retreat step. Since the set of admissible arcs is acyclic (by Lemma 10.6), after at most $2n$ advance steps of the first type, the algorithm will discover an admissible path and will perform an augmentation (because the longest path in the network has $2n$ nodes). Since the algorithm performs a total of $O(m \log U)$ augmentations, the number of advance steps of the first type is at most $O(nm \log U)$. The algorithm performs $O(nm)$ advance steps of the second type because each retreat step increases a node potential, and by Lemma 10.4, node potentials increase $O(n(n + m)) = O(nm)$ times. Therefore, the total number of advance steps is $O(nm \log U)$.

The amount of time needed to relabel nodes in N_1 is $O(n \sum_{i \in N} |A(i)|) = O(nm)$. The time needed to relabel nodes in N_2 is also $O(nm)$ since $|N_2| = m$ and the degree of each node in N_2 is constant (i.e., it is 2). The same arguments show that the algorithm requires $O(nm)$ time to identify admissible arcs. We have, therefore, established the following result.

Theorem 10.10. *The double scaling algorithm solves the minimum cost flow problem in $O(nm \log U \log(nC))$ time.* ♦

One nice feature of the double scaling algorithm is that it achieves an excellent worst-case running time for solving the minimum cost flow problem and yet is fairly simple, both conceptually and computationally.

10.5 MINIMUM MEAN CYCLE-CANCELING ALGORITHM

The three minimum cost flow algorithms we have discussed in this chapter—the capacity scaling algorithm, the cost scaling algorithm, and the double scaling algorithm—are weakly polynomial-time algorithms because their running times depend on $\log U$ and/or $\log C$. Although these algorithms are capable of solving any problem with integer or rational data, they are not applicable to problems with irrational data. In contrast, the running times of strongly polynomial-time algorithms depend only on n and m ; consequently, these algorithms are capable of solving problems with irrational data, assuming that a computer can perform additions and subtractions on irrational data. In this and the next two sections, we discuss several strongly polynomial time algorithms for solving any class of minimum cost flow problems, including those with irrational data.

The algorithm discussed in this section is a special case of the cycle-canceling algorithm that we discussed in Section 9.6. Because this algorithm iteratively cancels cycles (i.e., augments flows along cycles) with the minimum mean cost in the residual network, it is known as the *(minimum) mean cycle-canceling algorithm*. Recall from Section 5.7 that the *mean cost* of a directed cycle W is $(\sum_{(i,j) \in W} c_{ij}) / |W|$, and that the *minimum mean cycle* is a cycle with the smallest mean cost in

the network. In Section 5.7 we showed how to use dynamic programming algorithm to find the minimum mean cycle in $O(nm)$ time.

The minimum mean cycle-canceling algorithm starts with a feasible flow x in the network. At every iteration, the algorithm identifies a minimum mean cycle W in $G(x)$. If the mean cost of the cycle W is negative, the algorithm augments the maximum possible flow along W , updates $G(x)$, and repeats this process. If the mean cost of W is nonnegative, $G(x)$ contains no negative cycle and x is a minimum cost flow, so the algorithm terminates. This algorithm is surprisingly simple to state; even more surprisingly, the algorithm runs in strongly polynomial time.

To establish the worst-case complexity of the minimum mean cycle-canceling algorithm, we recall a few facts. In our subsequent discussion, we often use Property 9.2(b), which states that for any set of node potentials π and any directed cycle W , the sum of the costs of the arcs in W equals the sum of the reduced costs of the arcs in W . We will also use the following property concerning sequences of real numbers, which is a variant of the geometric improvement argument (see Section 3.3).

Property 10.11. *Let α be a positive integer and let y_1, y_2, y_3, \dots be a sequence of real numbers satisfying the condition $y_{k+1} \leq (1 - 1/\alpha)y_k$ for every k . Then for every value of k , $y_{k+\alpha} \leq y_k/2$.*

Proof. We first rewrite the expression $y_{k+1} \leq (1 - 1/\alpha)y_k$ as $y_k \geq y_{k+1} + y_{k+1}/(\alpha - 1)$. We now use this last expression repeatedly to replace the first term on the right-hand side, giving

$$\begin{aligned} y_k &\geq y_{k+1} + y_{k+1}/(\alpha - 1) \geq y_{k+2} + y_{k+2}/(\alpha - 1) + y_{k+1}/(\alpha - 1) \\ &\geq y_{k+2} + 2y_{k+2}/(\alpha - 1) \geq y_{k+3} + 3y_{k+3}/(\alpha - 1) \\ &\vdots \\ &\geq y_{k+\alpha} + \alpha y_{k+\alpha}/(\alpha - 1) \geq 2y_{k+\alpha}, \end{aligned}$$

which is the assertion of the property. \blacklozenge

We divide the worst-case analysis of the minimum mean cycle algorithm into two parts: First, we show that the algorithm is weakly polynomial-time; then we establish its strong polynomiality. Although the description of the algorithm does not use scaling techniques, the worst-case analysis borrows ideas from the cost scaling algorithm that we discussed in Section 10.3. In particular, the notion of ϵ -optimality discussed in that section plays a crucial role in its analysis. We will show that the flows maintained by the minimum mean cycle-canceling algorithm are ϵ -optimal flows satisfying the conditions that (1) between any two consecutive iterations the value of ϵ either stays the same or decreases; (2) occasionally, the value of ϵ strictly decreases; and (3) eventually, $\epsilon < 1/n$ and the algorithm terminates (see Lemma 10.2). As we observed in Section 10.3, the cost scaling algorithm's explicit strategy is to reduce ϵ from iteration to iteration. Although the minimum mean cycle-canceling algorithm also reduces the value of ϵ (although periodically, rather than at every iteration), the reduction is very much an implicit by-product of the algorithm.

We first establish a connection between the ϵ -optimality of a flow x and the

mean cost of a minimum mean cycle in $G(x)$. Recall that a flow x is ϵ -optimal if for some set of node potentials, the reduced cost of every arc is at least $-\epsilon$. Notice that any flow x will be ϵ -optimal for many values of ϵ , because a flow that is ϵ -optimal is also ϵ' -optimal for all $\epsilon' \geq \epsilon$. For any particular set of node potentials π , we let $\epsilon^\pi(x)$ be the negative of the minimum value of any reduced cost [i.e., $\epsilon^\pi(x) = -\min[c_{ij}^\pi : (i, j) \text{ in } G(x)]$]. Thus $c_{ij}^\pi \geq -\epsilon^\pi(x)$ and $c_{ij}^\pi = -\epsilon^\pi(x)$ for some arc (i, j) . Thus x is ϵ -optimal for $\epsilon = \epsilon^\pi(x)$. Potentially, we could find a smaller value of ϵ by using other values of the node potentials. With this thought in mind, we let $\epsilon(x) = \min_\pi \epsilon^\pi(x)$. Note that $\epsilon(x)$ is the smallest value of ϵ for which the flow x is ϵ -optimal. As additional notation, we let $\mu(x)$ denote the mean cost of the minimum mean cycle in $G(x)$.

Note that since x is $\epsilon(x)$ -optimal, conditions (10.2) imply that $\sum_{(i,j) \in W} c_{ij} = \sum_{(i,j) \in W} c_{ij}^\pi \geq -\epsilon^\pi(x) |W|$. Choosing W as the minimum mean cycle and dividing this expression by $|W|$, we see that $\mu(x) \geq -\epsilon(x)$. As we have seen, this inequality is a simple consequence of the definitions of ϵ -optimality and of the minimum mean cycle cost; it uses the fact that if we can bound the reduced cost of every arc around a cycle, this same bound applies to the average cost around the cycle. Perhaps surprisingly, however, we can obtain a converse result: that is, we can always find a set of node potentials so that every arc around the minimum mean cycle has the same reduced cost and that this cost equals $-\epsilon(x)$. Our next two results establish this property.

Lemma 10.12. *Let x be a nonoptimal flow. Then $\epsilon(x) = -\mu(x)$.*

Proof. Since our observation in the preceding paragraph shows that $\epsilon(x) \geq -\mu(x)$, we only need to show that $\epsilon(x) \leq -\mu(x)$.

Let W be a minimum mean cycle in the residual network $G(x)$, and let $\mu(x)$ be the mean cost of this cycle. Suppose that we replace each arc cost c_{ij} by $c'_{ij} = c_{ij} - \mu(x)$. This transformation reduces the mean cost of every directed cycle in $G(x)$ by $\mu(x)$ units. Consequently, the minimum mean cost of the cycle W becomes zero, which implies that the residual network contains no negative cost cycle. Let $d'(\cdot)$ denote the shortest path distances in $G(x)$ from a specified node s to all other nodes with c'_{ij} as the arc lengths. The shortest path optimality conditions imply that

$$d'(j) \leq d'(i) + c'_{ij} = d'(i) + c_{ij} - \mu(x) \quad \text{for each arc } (i, j) \text{ in } G(x). \quad (10.7)$$

If we let $\pi(j) = d'(j)$, then (10.7) becomes

$$c_{ij}^\pi \geq \mu(x) \quad \text{for each arc } (i, j) \text{ in } G(x), \quad (10.8)$$

which implies that x is $(-\mu(x))$ -optimal. Therefore, $\epsilon(x) \leq -\mu(x)$, completing the proof of the lemma. \blacklozenge

Lemma 10.13. *Let x be any nonoptimal flow. Then for some set of node potentials π , $c_{ij}^\pi = \mu(x) = -\epsilon(x)$ for every arc (i, j) in the minimum mean cycle W of $G(x)$.*

Proof. Let π be defined as in the proof of the preceding lemma; with these set of node potentials, the reduced costs satisfy (10.8). The cost of the cycle W equals $\sum_{(i,j) \in W} c_{ij}$, which also equals its reduced cost $\sum_{(i,j) \in W} c_{ij}^\pi$. Con-

sequently, $\sum_{(i,j) \in W} c_{ij}^{\pi} = \mu(x) |W|$. This equation and (10.8) imply that $c_{ij}^{\pi} = \mu(x)$ for each arc (i, j) in W . Lemma 10.12 establishes that $c_{ij}^{\pi} = -\epsilon(x)$ for every arc in W . ♦

We next show that during the execution of the minimum mean cycle-canceling algorithm, $\epsilon(x)$ never increases; moreover, within m consecutive iterations $\epsilon(x)$ decreases by a factor of at least $(1 - 1/n)$.

Lemma 10.14. *For a nonoptimal flow x , if we cancel a minimum mean cycle in $G(x)$, $\epsilon(x)$ cannot increase [alternatively, $\mu(x)$ cannot decrease].*

Proof. Let W denote the minimum mean cycle in $G(x)$. Lemma 10.13 implies that for some set of node potentials π , $c_{ij}^{\pi} = -\epsilon(x)$ for each arc $(i, j) \in W$. Let x' denote the flow obtained after we have canceled the cycle W . This flow augmentation deletes some arcs in W from the residual network and adds some other arcs, which are reversals of the arcs in W . Consider any arc (i, j) in $G(x')$. If (i, j) is in $G(x)$, then, by hypothesis, $c_{ij}^{\pi} \geq -\epsilon(x)$. If (i, j) is not in $G(x)$, then (i, j) is a reversal of some arc (j, i) in $G(x)$ for which $c_{ji}^{\pi} = -\epsilon(x)$. Therefore, $c_{ij}^{\pi} = -c_{ji}^{\pi} = \epsilon(x) > 0$. In either case, $c_{ij}^{\pi} \geq -\epsilon(x)$ for each arc (i, j) in $G(x')$. Consequently, the minimum mean cost of any cycle in $G(x')$ will be at least $-\epsilon(x)$, since the mean cost around a cycle, which equals the mean reduced cost, must be at least as large as the minimum value of the reduced costs. Therefore, in light of Lemma 10.12, as asserted, $\epsilon(x') = \mu(x') \geq -\epsilon(x) = \mu(x)$. ♦

Lemma 10.15. *After a sequence of m minimum mean cycle cancellations starting with a flow x , the value of the optimality parameter $\epsilon(x)$ decreases to a value at most $(1 - 1/n)^m \epsilon(x)$ [i.e., to at most $(1 - 1/n)^m$ times its original value].*

Proof. Let π denote a set of node potentials satisfying the conditions $c_{ij}^{\pi} \geq -\epsilon(x)$ for each arc (i, j) in $G(x)$. For convenience, we designate those arcs in $G(x)$ with (strictly) negative reduced costs as *negative arcs* (with respect to the reduced costs). We now classify the subsequent cycle cancellations into two types: (1) all the arcs in the canceled cycle are negative (a type 1 cancellation), and (2) at least one arc in the canceled cycle has a nonnegative reduced cost (a type 2 cancellation). We claim that the algorithm will perform at most m type 1 cancellations before it either terminates or performs a type 2 cancellation. This claim follows from the observations that each type 1 cancellation deletes at least one negative arc from the (current) residual network and all the arcs that the cancellation adds to the residual network have positive reduced cost with respect to π (as shown in the proof of Lemma 10.14). Consequently, if within m iterations, the algorithm performs no type 2 cancellations, all the arcs in the residual network will have nonnegative reduced costs with respect to π and the algorithm will terminate with an optimal flow.

Now consider the first time the algorithm performs a type 2 cancellation. Suppose that the algorithm cancels the cycle W , which contains at least one arc with a nonnegative reduced cost; let x' and x'' denote the flows just before and after the cancellation. Then $c_{ij}^{\pi} \geq -\epsilon(x')$ for each arc $(i, j) \in W$ and $c_{kl}^{\pi} \geq 0$ for some arc $(k, l) \in W$. As a result, since $c(W) = \sum_{(i,j) \in W} c_{ij}^{\pi}$, the cost $c(W)$ of W with respect to the flow x' satisfies the condition $c(W) \geq [(|W| - 1)(-\epsilon(x'))]$. By Lemma 10.14,

the cancelation cannot increase the minimum mean cost and, therefore, $\mu(x'') \geq \mu(x')$. But since $\mu(x')$ is the mean cost of W with respect to x' , $\mu(x'') \geq \mu(x') \geq (1 - 1/n) \sum_{i \in W} (-\epsilon(x'_i)) \geq (1 - 1/n) \sum_{i \in W} (-\epsilon(x''_i))$. This inequality implies that $-\mu(x'') \leq (1 - 1/n) \sum_{i \in W} \epsilon(x''_i)$. Using the fact that $\mu(x'') = -\epsilon(x'')$, we see that $\epsilon(x'') \leq (1 - 1/n) \epsilon(x')$. This result establishes the lemma. ♦

As indicated by the next theorem, the preceding two lemmas imply that the minimum mean cycle-canceling algorithm performs a polynomial number of iterations.

Theorem 10.16. *If all arc costs are integer, the minimum mean cycle-canceling algorithm performs $O(nm \log(nC))$ iterations and runs in $O(n^2 m^2 \log(nC))$ time.*

Proof. Let x denote the flow at any point during the execution of the algorithm. Initially, $\epsilon(x) \leq C$ because every flow is C -optimal (see Lemma 10.2). In every m consecutive iterations, the algorithm decreases $\epsilon(x)$ by a factor of $(1 - 1/n)$. When $\epsilon(x) < 1/n$, the algorithm terminates with an optimal flow (see Lemma 10.2). Therefore, the algorithm needs to decrease $\epsilon(x)$ by a factor of nC over all iterations. By Lemma 10.15, the mean cost of a cycle becomes smaller by a factor of at least $(1 - 1/n)$ in every m iterations. Property 10.11 implies that the minimum mean cycle cost decreases by a factor of 2 every nm iterations, so that within $nm \log(nC)$ iterations, the minimum mean cycle cost decreases from C to $1/n$. At this point the algorithm terminates with an optimal flow. This conclusion establishes the first part of the theorem. Since the bottleneck operation in each iteration is identifying a minimum mean cycle, which requires $O(nm)$ time (see Section 5.7), we also have established the second part of the theorem.

Having proved that the minimum mean cycle-canceling algorithm runs in polynomial time, we next obtain a strongly polynomial bound on the number of iterations the algorithm performs. Our analysis rests upon the following rather useful result: If the absolute value of the reduced cost of an arc (k, l) is “significantly greater than” the current value of the parameter $\epsilon(x)$, the flow on the arc (k, l) in any optimal solution is the same as the current flow on this arc. In other words, the flow on the arc (k, l) becomes “fixed.” As we will show, in every $O(nm \log n)$ iterations, the algorithm will fix at least one additional arc at its lower bound or at its upper bound. As a result, within $O(nm^2 \log n)$ iterations, the algorithm will have fixed all the arcs and will terminate with an optimal flow.

We define an arc to be ϵ -fixed if the flow on this arc is the same for all ϵ' -optimal flows whenever $\epsilon' \leq \epsilon$. Since the value of $\epsilon(x)$ of the $\epsilon(x)$ -optimal flows, that the minimum mean cycle-canceling algorithm maintains, is nonincreasing, the flow on an $\epsilon(x)$ -fixed arc will not change during the execution of the algorithm and will be the same in every optimal flow. We next establish a condition that will permit us to fix an arc.

Lemma 10.17. *Suppose that x is an $\epsilon(x)$ -optimal flow with respect to the potentials π , and suppose that for some arc $(k, l) \in A$, $|c_{kl}^\pi| \geq 2n\epsilon(x)$. Then arc (k, l) is an $\epsilon(x)$ -fixed arc.*

Proof. Let $\epsilon = \epsilon(x)$. We first prove the lemma when $c_{kl}^{\pi} \geq 2n\epsilon$. The ϵ -optimality condition (10.1a) implies that $x_{kl} = 0$. Suppose that some $\epsilon(x')$ -optimal flow x' , with $\epsilon(x') \leq \epsilon(x)$, satisfies the condition that $x'_{kl} > 0$. The flow decomposition theorem (i.e., Theorem 3.5) implies that we can express x' as x plus the flow along at most m augmenting cycles in $G(x)$. Since $x_{kl} = 0$ and $x'_{kl} > 0$, one of these cycles, say W , must contain the arc (k, l) as a forward arc. Since each arc $(i, j) \in W$ is in the residual network $G(x)$, and so satisfies the condition $c_{ij}^{\pi} \geq -\epsilon$, the reduced cost (or, cost) of the cycle W is at least $c_{kl}^{\pi} - \epsilon(|W| - 1) \geq 2n\epsilon - \epsilon(n - 1) > n\epsilon$.

Now consider the cycle W' obtained by reversing the arcs in W . The cycle W' must be a directed cycle in the residual network $G(x')$ (see Exercise 10.6). The cost of the cycle W' is the negative of the cost of the cycle W and so must be less than $-n\epsilon \leq -n\epsilon(x')$. Therefore, the mean cost of W' is less than $-\epsilon(x')$. Lemma 10.12 implies that x' is not $\epsilon(x')$ -optimal, which is a contradiction.

We next consider the case when $c_{kl}^{\pi} \leq -2n\epsilon$. In this case the ϵ -optimality condition (10.1c) implies that $x_{kl} = u_{kl}$. Using an analysis similar to the one used in the preceding case, we can show that no ϵ -optimal flow x' can satisfy the condition $x'_{kl} < u_{kl}$. ♦

We are now in a position to obtain a strongly polynomial bound on the number of iterations performed by the minimum mean cycle-canceling algorithm.

Theorem 10.18. *For arbitrary real-valued arc costs, the minimum mean cycle-canceling algorithm performs $O(nm^2 \log n)$ iterations and runs in $O(n^2 m^3 \log n)$ time.*

Proof. Let $K = nm(\lceil \log n \rceil + 1)$. We divide the iterations performed by the algorithm into groups of K consecutive iterations. We claim that each group of iterations fixes the flow on an additional arc (k, l) (i.e., the iterations after those in the group do not change the value of x_{kl}). The theorem follows immediately from this claim, since the algorithm can fix at most m arcs, and each iteration requires $O(nm)$ time.

Consider any group of iterations. Let x be the flow before the first iteration of the group and let x' be the flow after the last iteration of the group. Let $\epsilon = \epsilon(x)$, $\epsilon' = \epsilon(x')$, and let π' be the node potentials for which x' satisfies the ϵ' -optimality conditions. Since every nm iterations reduce ϵ by a factor of at least 2, the $nm(\lceil \log n \rceil + 1)$ iterations between x and x' reduce ϵ by a factor of at least $2^{\lceil \log n \rceil + 1}$. Therefore, $\epsilon' \leq (\epsilon/2^{\lceil \log n \rceil + 1}) \leq \epsilon/2n$. Alternatively, $-\epsilon \leq -2n\epsilon'$.

Let W be the cycle canceled when the flow has value x . Lemma 10.12 and the fact that the sum of the costs and reduced costs around every cycle are the same, imply that for any values of the node potentials, the average reduced cost around the cycle W equals $\mu(x) = -\epsilon$. Therefore, with respect to the potentials π' , at least one arc (k, l) in W must have a reduced cost as small as $-\epsilon$, so $c_{kl}^{\pi'} = -\epsilon \leq -2n\epsilon'$ for some arc (k, l) in W . By Lemma 10.17, the flow on arc (k, l) will not change in any subsequent iteration. Next notice that in the first iteration in the group, the algorithm changed the value of x_{kl} . Thus each group fixes the flow on at least one additional arc, completing the proof of the theorem. ♦

We might conclude this section with a few observations. First, note that we need not formally compute the value of $\epsilon(x)$ at each iteration, nor do we need to identify the ϵ -fixed arcs at any stage in the algorithm. Indeed, we can use any method to find the minimum mean cost cycle at each step; in principle, we need not maintain or ever compute any reduced costs. As we noted earlier in this section, the minimum mean cycle-canceling algorithm implicitly reduces $\epsilon(x)$ and fixes some arcs as it proceeds—we need not keep track of the algorithm’s progress concerning these features.

We also might note that the ideas presented in this section would also permit us to develop a strongly polynomial-time version of the cost scaling algorithm that we discussed in Section 10.3. In Exercise 10.12 we consider this modification of the cost scaling algorithm and analyze its running time.

10.6 REPEATED CAPACITY SCALING ALGORITHM

The minimum cost flow problem described in Section 10.5 uses the idea that whenever the reduced cost of an arc is *sufficiently large*, we can “fix” the flow on the arc. By incorporating a similar idea in the capacity scaling algorithm, we can develop another strongly polynomial time algorithm. As we will see, when the flow on an arc (i, j) is *sufficiently large*, the potentials of nodes i and j become “fixed” with respect to each other. In this section we discuss the details of this algorithm, which we call the *repeated capacity scaling algorithm*.

The repeated capacity scaling algorithm to be discussed in this section is different from all the other minimum cost flow algorithms discussed in this book. All of the other algorithms solve the primal minimum cost flow problem (9.1) and obtain an optimal flow; the repeated capacity scaling algorithm solves the dual minimum cost flow problem (9.10). This algorithm obtains an optimal set of node potentials for (9.10) and then uses it to determine an optimal flow.

The repeated capacity scaling algorithm is a modified version of the capacity scaling algorithm discussed in Section 10.2. For simplicity, we describe the algorithm for the uncapacitated minimum cost flow problem; we could solve the capacitated problem by converting it to the uncapacitated problem using the transformation described in Section 2.4. Recall that in the capacity scaling algorithm, each arc flow is an integral multiple of the scale factor Δ . For uncapacitated networks, each residual capacity r_{ij} is also an integral multiple of Δ , because either $r_{ij} = u_{ij} = \infty$, or $r_{ij} = x_{ji} = k\Delta$ for some integer k . This observation implies that the Δ -residual network $G(x, \Delta)$ is the same as the residual network $G(x)$. As a result, the algorithm for the uncapacitated problem does not require the preprocessing (i.e., saturating the arcs violating the optimality conditions) at the beginning of each scaling phase. The following property is an immediate consequence of this result.

Property 10.19. *The capacity scaling algorithm for the uncapacitated minimum cost flow problem satisfies the following properties: (a) the excesses at the nodes are monotonically decreasing; (b) the sum of the excesses at the beginning of the Δ -scaling phase is at most $2n\Delta$; and (c) the algorithm performs at most $2n$ augmentations per scaling phase.*

The repeated capacity scaling algorithm is based on the three simple results stated in the following lemmas.

Lemma 10.20. *Suppose that at the beginning of the Δ -scaling phase, $b(k) > 6n^2\Delta$ for some node $k \in N$. Then some arc (k, l) with $x_{kl} > 4n\Delta$ emanates from node k .*

Proof. Property 10.19 implies that at the beginning of the Δ -scaling phase, the sum of the excesses is at most $2n\Delta$. Therefore, $e(k) \leq 2n\Delta$. Since $b(k) > 6n^2\Delta$ and $e(k) \leq 2n\Delta$, the net outflow of node k [i.e., $b(k) - e(k)$] is strictly greater than $(6n^2\Delta - 2n\Delta)$. Since fewer than n arcs emanate from node k , the flow on at least one of these arcs must be strictly more than $(6n^2\Delta - 2n\Delta)/n \geq (4n^2\Delta)/n = 4n\Delta$, which concludes the lemma. ♦

Lemma 10.21. *If at the beginning of the Δ -scaling phase $x_{kl} > 4n\Delta$, then for some optimal solution $x_{kl} > 0$.*

Proof. Property 10.19 implies that the algorithm performs at most $2n$ augmentations in each scaling phase. The fact that the algorithm augments exactly Δ units of flow in every augmentation in the Δ -scaling phase implies that the total flow change due to all augmentations in the subsequent scaling phases is at most $2n(\Delta + \Delta/2 + \Delta/4 + \cdots + 1) < 4n\Delta$. Consequently, if $x_{kl} > 4n\Delta$ at the beginning of the Δ -scaling phase, then $x_{kl} > 0$ when the algorithm terminates. ♦

Lemma 10.22. *Suppose that $x_{kl} > 0$ in an optimal solution of the minimum cost flow problem. Then with respect to every set of optimal node potentials, the reduced cost of arc (k, l) is zero.*

Proof. Suppose that x satisfies the complementary slackness optimality condition (9.8) with respect to the node potential π . The condition (9.8b) implies that $c_{kl}^\pi = 0$. Property 9.8 implies that if x satisfies the complementary slackness optimality condition (9.8b) with respect to some node potential, it satisfies this condition with respect to every optimal node potential. Consequently, the reduced cost of arc (k, l) is zero with respect to every set of optimal node potentials. ♦

We are now in a position to discuss the essential ideas of the repeated capacity scaling algorithm. Let \mathbf{P} denote the minimum cost flow problem stated in (9.1). The algorithm applies the capacity scaling algorithm stated in Figure 10.1 to the problem \mathbf{P} . We will show that within $O(\log n)$ scaling phases, $b(k) > 6n^2\Delta$ for some node k and, by Lemma 10.20, some arc (k, l) satisfies the condition $x_{kl} > 4n\Delta$. Lemmas 10.21 and 10.22 imply that for any set of optimal node potentials, the reduced cost of arc (k, l) will be zero. This result allows us to show, as described next, that we can *contract* the nodes k and l into a single node, thereby obtaining a new minimum cost flow problem defined on a network with one fewer node.

Suppose that we are using the capacity scaling algorithm to solve a minimum cost flow problem \mathbf{P} with arc costs c_{ij} and at some stage we realize that for an arc (k, l) , $x_{kl} > 4n\Delta$. Let π denote the node potentials at this point. The optimality condition (9.8b) implies that

$$c_{kl} - \pi(k) + \pi(l) = 0. \quad (10.9)$$

Now consider the same minimum cost flow problem, but with the cost of each arc (i, j) equal to $c'_{ij} = c_{ij} - \pi(i) + \pi(j)$. Let \mathbf{P}' denote the modified minimum cost flow problem. Condition (10.9) implies that

$$c'_{kl} = 0. \quad (10.10)$$

We next observe that the problems \mathbf{P} and \mathbf{P}' have the same optimal solutions (see Property 2.4 in Section 2.4). Since $x_{kl} > 4n\Delta$, Lemmas 10.21 and 10.22 imply that in problem \mathbf{P}' the reduced cost of arc (k, l) will be zero. If π' denotes an optimal set of node potentials for \mathbf{P}' , then

$$c'_{kl} - \pi'(k) + \pi'(l) = 0. \quad (10.11)$$

Substituting (10.10) in (10.11) implies that $\pi'(k) = \pi'(l)$.

The preceding discussion shows that if $x_{kl} > 4n\Delta$ for some arc (k, l) , we can “fix” one node potential with respect to the other. The discussion also shows that if we solve the problem \mathbf{P}' with the additional constraint that the potentials of nodes k and l are same, this constraint will not eliminate the optimal solution of \mathbf{P}' . But how can we solve a minimum cost flow problem when two node potentials must be the same?

Consider the dual minimum cost flow problem stated in (9.10). In this problem we replace both $\pi(k)$ and $\pi(l)$ by $\pi(p)$. This substitution gives us a linear programming problem with one less dual variable (or, node potential). The reader can easily verify that the resulting problem is a dual minimum cost flow problem on the network with nodes k and l contracted into a single node p . The contraction operation consists of (1) letting $b(p) = b(k) + b(l)$, (2) replacing each arc (i, k) or (i, l) by the arc (i, p) , (3) replacing each arc (k, i) or (l, i) by the arc (p, i) , and (4) letting the cost of an arc in the contracted network equal that of the arc it replaces. We point out that the contraction might produce multiarcs (i.e., more than one arc with the same tail and head nodes). The purpose of contraction operations should be clear; since each contraction operation reduces the size of the network by one node, we can apply at most n of these operations.

We can now describe the repeated capacity scaling algorithm. We first compute $U = \max\{b(i) : i \in N \text{ and } b(i) > 0\}$ and initialize $\Delta = 2^{\lceil \log U \rceil}$. Let node k be a node with $b(k) = U$. We then apply the capacity scaling algorithm as described in Figure 10.1. Each scaling phase of the capacity scaling algorithm decreases Δ by a factor of 2; therefore, since the initial value of Δ is $b(k)$, after at most $q = \log(6n^2) = O(\log n)$ phases, $\Delta = b(k)/2^q \leq b(k)/6n^2$. The algorithm might obtain a feasible flow before $\Delta \leq b(k)/6n^2$ (in which case it terminates); if not, then by Lemma 10.20, some arc (k, l) will satisfy the condition that $x_{kl} > 4n\Delta$. The algorithm then defines a new minimum cost flow problem with nodes k and l contracted into a new node p , and the cost of each arc is the reduced cost of the corresponding arc before the contraction. We solve the new minimum cost flow problem afresh by redefining U as the largest supply in the contracted network and reapplying the capacity scaling algorithm described in Figure 10.1. We repeat these steps until the algorithm terminates. The algorithm terminates in one of the two ways: (1) while applying the capacity scaling algorithm, it obtains a flow; or (2) it contracts the network into a

single node p [with $b(p) = 0$], which is trivially solvable by a zero flow. At this point we expand the contracted nodes and obtain an optimal flow in the expanded network. We show how to expand the contracted nodes a little later. The preceding discussion shows that the algorithm performs $O(n \log n)$ scaling phases, and since each scaling phase solves at most $2n$ shortest path problems, the running time of the algorithm is $O(n^2 \log n S(n, m))$. In this expression, $S(n, m)$ is the minimum time required by a strongly polynomial-time algorithm for solving a shortest path problem with nonnegative arc lengths. [Recall from Chapter 4 that $O(m + n \log n)$ is currently the best known such bound.]

We illustrate the repeated capacity scaling algorithm on the example shown in Figure 10.11(a). When applied to this example, the capacity scaling algorithm performs 100 scaling phases with $\Delta = 2^{99}, 2^{98-1}, \dots, 2^0$. The strongly polynomial version, however, terminates within five phases, as shown next.

Phase 1. In this phase, $\Delta = 2^{99}$, $S(\Delta) = \{1, 2\}$, and $T(\Delta) = \{3, 4\}$. The algorithm augments Δ units of flow along the two paths 1–3 and 2–1–3–4. Figure 10.11(b) shows the solution at the end of this phase.

Phase 2. In this phase, $\Delta = 2^{98}$. The algorithm augments Δ units of flow along the path 1–3.

Phase 3. In this phase, $\Delta = 2^{97}$. The algorithm augments Δ units of flow along the path 1–3.

Phase 4. In this phase, $\Delta = 2^{96}$. The algorithm finds that the flow on the arc $(1, 3)$ is $2^{100} + 2^{99} + 2^{98}$, which is more than $4n\Delta = 2^{100}$. Therefore, the algorithm contracts the nodes 1 and 3 into a new node 5 and obtains the minimum cost flow problem shown in Figure 10.11(c), which it then proceeds to solve.

Phase 5. In this phase, $\Delta = 2^{95}$. The algorithm augments Δ units of flow along the path 2–5–4. The solution is a flow now; consequently, the algorithm terminates. The corresponding flow in the original network is $x_{21} = 2^{99}$, $x_{13} = 2^{100} - 1$, and $x_{34} = 2^{99}$.

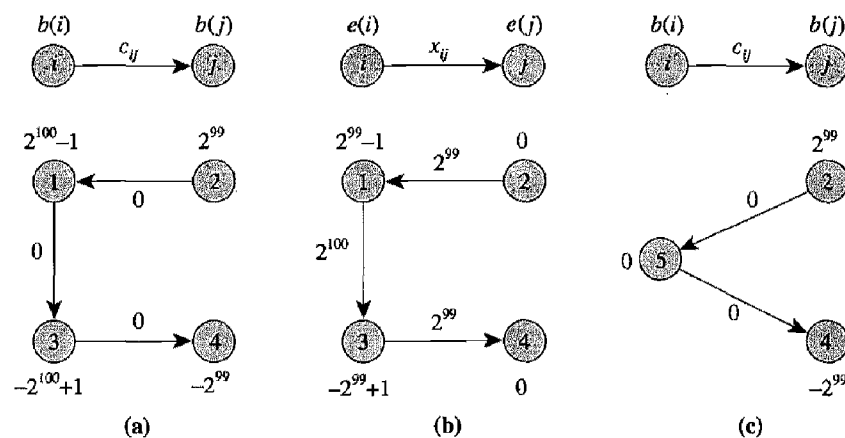


Figure 10.11 Illustrating the repeated capacity scaling algorithm: (a) minimum cost flow problem; (b) solution after the first phase; (c) minimum cost flow problem after contracting the nodes 1 and 3 into a new node 5.

We now explain how we expand the contracted network, and in the process we prove that the algorithm determines an optimal solution of the minimum cost flow problem. The algorithm, in fact, first determines an optimal set of node potentials of the problem, and then by solving a maximum flow problem (as described in Section 9.5) determines an optimal flow. The algorithm obtains an optimal set of node potentials for the original problem by repeated use of the following result.

Property 10.23. *Let \mathbf{P} be a problem with arc costs c_{ij} and \mathbf{P}' be the same problem with arc costs $c_{ij} - \pi(i) + \pi(j)$. If π' is an optimal set of node potentials for problem \mathbf{P}' , then $\pi + \pi'$ is an optimal set of node potentials for \mathbf{P} .*

Proof. This property easily follows from the observation that if a solution x satisfies the reduced cost optimality condition (9.7) with respect to the arc costs $c_{ij} - \pi(i) + \pi(j)$ and node potentials π' , the same solution satisfies these conditions with arc costs c_{ij} and node potentials $\pi + \pi'$. ♦

We expand (or uncontract) the nodes in the reverse order in which we contracted them in the strongly polynomial algorithm and obtain optimal node potentials of the successive problems. In earlier stages, between two successive problems, we performed two transformations in the following order: (1) we replaced the arc cost c_{ij} by its reduced cost $c_{ij} - \pi(i) + \pi(j)$, and (2) we contracted two nodes k and l into a single new node p . We undo these transformations in the reverse order. To undo the contracted node p , for case (2) we set the potentials of nodes k and l equal to that of node p , and for case (1) we add π to the existing node potentials. When we have expanded all the contracted nodes, the resulting node potentials are an optimal set of node potentials for the minimum cost flow problem. Then, as described in Section 9.5, we can use these node potentials to obtain an optimal flow by solving a maximum flow problem. The following theorem summarizes the preceding discussion.

Theorem 10.24. *The repeated capacity scaling algorithm solves the uncapacitated minimum cost flow problem in $O(n^2 \log n S(n, m))$ time.* ♦

Since the best known strongly polynomial-time algorithm for solving the shortest path problem with nonnegative arc lengths runs in $O(m + n \log n)$ time, the best current bound for the uncapacitated minimum cost flow problem is $O(n \log n(m + n \log n))$. We can solve the capacitated minimum cost flow problem by the repeated capacity scaling algorithm by first transforming it to an uncapacitated problem (see Section 2.4). The uncapacitated network will have $n' = n + m$ nodes and $m' = 2m$ arcs. When applied to this network, the repeated capacity scaling algorithm will perform $O(n' \log n') = O(m \log n)$ scaling phases and solve $O(m') = O(m)$ shortest path problems in each scaling phase. Thus the running time of the algorithm is the time needed to solve $O(m^2 \log n)$ shortest path problems. Each shortest path problem in the uncapacitated network requires $O(2m + (m + n) \log (m + n)) = O(m + m \log n)$ time, but using a clever approach for solving the resulting shortest path problem (as discussed in Exercise 4.53) we can obtain a better bound of $O(m + n \log n)$. Consequently, the repeated capacity scaling algorithm requires $O(m^2 \log n(m + n \log n))$ time to solve a capacitated minimum cost flow problem.

10.7 ENHANCED CAPACITY SCALING ALGORITHM

In this section we discuss yet another strongly polynomial-time algorithm for the minimum cost flow problem. This algorithm is a variant of the capacity scaling algorithm that we discussed in Section 10.2 and draws on some ideas from the repeated capacity scaling algorithm discussed in Section 10.6. We refer to this algorithm as the *enhanced capacity scaling algorithm*. This algorithm runs in $O((m \log n)(m + \log n))$ time for the capacitated minimum cost flow problem and is currently the fastest strongly polynomial-time algorithm for solving the minimum cost flow problem. In this section we first show how to solve the enhanced capacity scaling algorithm for the *uncapacitated* minimum cost flow problem; we can solve the capacitated problem by transforming it to an uncapacitated problem (see Section 2.4).

Recall from Section 10.6 that the essential idea in the repeated capacity scaling algorithm is to identify arcs with *sufficiently large* flow. The repeated capacity scaling algorithm identifies such an arc (k, l) within $O(\log n)$ scaling phases, contracts the nodes k and l into a single node, and solves the resulting minimum cost flow problem afresh. For the uncapacitated minimum cost flow problem, this algorithm performs a total of $O(n \log n)$ scaling phases and $O(n^2 \log n)$ shortest path augmentations. The enhanced capacity scaling algorithm adopts a similar approach but it differs in the following two ways: (1) the algorithm does not explicitly perform the contraction operation; and (2) the algorithm does not solve the minimum cost flow problem afresh, but continues from where it left off in its earlier computations. By avoiding contractions, the algorithm achieves ease of coding (because contractions change the network structure and so its computer representation) and maintains a pseudoflow satisfying the dual optimality conditions at every step until the end, at which point it becomes an optimal flow. Moreover, the total number of scaling phases is $O(n \log n)$ and the total number of shortest path augmentations in these scalings phases is also $O(n \log n)$. Consequently, if $S(n, m)$ is the time required to solve a shortest path problem with nonnegative arc lengths, the running time of the enhanced capacity scaling algorithm for uncapacitated problems is $O(n \log n S(n, m))$. For capacitated minimum cost flow problems, this time bound becomes $O(m \log n S(n, m))$. [By Exercise 4.53 the time bound for the shortest path problem in the transformed network is $O(S(n, m))$ rather than $O(S(n + m, 2m))$ even though the transformed network has $n + m$ nodes and $2m$ arcs.]

The enhanced capacity scaling algorithm proceeds by performing scaling phases for different values of the scale factor Δ . In the Δ -scaling phase, we say that an arc (i, j) has a *sufficiently large* flow if $x_{ij} \geq 8n\Delta$. [We later show that if $x_{ij} \geq 8n\Delta$, then arc (i, j) will have positive flow during the entire execution of the algorithm.] We refer to an arc with sufficiently large flow as an *abundant arc*; otherwise, we call it a *nonabundant arc*. We refer to the subgraph consisting of the node set N and abundant arcs as the *abundant subgraph*. The abundant subgraph typically contains several components, which we call *abundant components*. If the network contains no abundant arc, the abundant subgraph contains n components, each consisting of a singleton node. For simplicity, we will designate an abundant component by the set S of nodes it spans. We let $b(S) = \sum_{i \in S} b(i)$ and $e(S) = \sum_{i \in S} e(i)$.

We designate an (arbitrary) node in each abundant component as its *root* and refer to all the other nodes as *nonroot nodes*. By convention we assume that the

minimum index node in an abundant component is its root. For example, if $S = \{3, 5, 9\}$, then node 3 is the root node of the abundant component S . Throughout its execution, the enhanced capacity scaling algorithm satisfies the following properties.

Property 10.25 (Flow Property). *In the Δ -scaling phase, the flow on each non-abundant arc is an integral multiple of Δ ; an abundant arc can have any nonnegative flow value.*

Property 10.26 (Imbalance Property). *Each nonroot node has a zero imbalance; a root node can have an excess or a deficit.*

At the beginning of the enhanced capacity scaling algorithm, the network has no abundant arc and the abundant subgraph contains n components, each consisting of a singleton node. As the algorithm proceeds, it identifies abundant arcs and adds them to the abundant subgraph. Suppose that the algorithm adds a new abundant arc (i, j) at some stage. Let S_i and S_j , respectively, denote the abundant components containing the nodes i and j . If $S_i = S_j$ [i.e., the arc (i, j) has both of its endpoints in the same component], this addition does not create any new abundant component; otherwise, the addition creates a new abundant component consisting of the union of S_i and S_j . We refer to this operation as a *merge operation* because it merges the components S_i and S_j into a single abundant component. Notice that since each merge operation reduces the number of abundant components by one, the algorithm can perform at most n merge operations.

Whenever the algorithm merges the components S_i and S_j , we need to ensure that the solution satisfies the imbalance property. Suppose that i_r and j_r denote the root nodes of the components S_i and S_j before the merge operation. Suppose further that $i_r < j_r$. If $e(j_r) = 0$, after the merge operation the abundant subgraph satisfies the imbalance property. However, if $e(j_r)$ is nonzero, we satisfy the imbalance property by sending $e(j_r)$ units of flow from node j_r to node i_r using any path in the merged component. [Notice that if $e(j_r) < 0$, we should view this augmentation as augmenting $|e(j_r)|$ units of flow from node i_r to j_r so we eliminate the imbalance at node j_r .] Observe that this augmentation changes the flow on some abundant arcs by $|e(j_r)|$ units. We refer to this augmentation as an *imbalance-property augmentation*. In Exercise 10.26 we ask the reader to show how to perform merge operations and the subsequent imbalance-property augmentations in $O(m)$ time.

We are now in a position to describe the enhanced capacity scaling algorithm. Figure 10.12 gives an algorithmic description of this algorithm.

The enhanced capacity scaling algorithm performs two types of augmentations. The first type of augmentation enforces the imbalance property when the algorithm identifies new abundant arcs; we have earlier defined these augmentations as the *imbalance-property augmentations*. The second type of augmentation takes place from excess nodes to deficit nodes along shortest paths. We refer to these augmentations as *shortest-path augmentations*.

As we have already mentioned, the enhanced capacity scaling algorithm is a variant of the capacity scaling algorithm. These two algorithms differ in the following respects:

```

algorithm enhanced capacity scaling;
begin
  set  $x := 0$ ,  $\pi := 0$ , and  $e := b$ ;
  set  $\Delta := \max\{|e(i)| : i \in N\}$ ;
  while the residual network  $G(x)$  contains a node  $i$  with  $e(i) > 0$  do
    begin
      if  $\max\{e(i) : i \in N\} \leq \Delta/(8n)$  then  $\Delta := \max\{e(i) : i \in N\}$ ;
      {the  $\Delta$ -scaling phase begins here}
      for each nonabundant arc  $(i, j)$  do
        if  $x_{ij} \geq 8n\Delta$  then designate arc  $(i, j)$  as an abundant arc;
      update abundant components and reinstate the imbalance property;
      while the residual network  $G(x)$  contains a node  $k$  with  $|e(k)| \geq (n-1)\Delta/n$  do
        begin
          select a pair of nodes  $k$  and  $l$  satisfying the property that (i) either  $e(k) > (n-1)\Delta/n$ 
            and  $e(l) < -\Delta/n$ , or (ii)  $e(k) > \Delta/n$  and  $e(l) < -(n-1)\Delta/n$ ;
          considering reduced costs as arc lengths, compute shortest path distance  $d(\cdot)$  in
             $G(x)$  from node  $k$  to all other nodes;
           $\pi(i) := \pi(i) - d(i)$  for all  $i \in N$ ;
          augment  $\Delta$  units of flow along the shortest path in  $G(x)$  from node  $k$  to node  $l$ ;
        end;
      {the  $\Delta$ -scaling phase ends here}
       $\Delta := \Delta/2$ ;
    end;
  end;

```

Figure 10.12 Enhanced capacity scaling algorithm.

1. In the capacity scaling algorithm, we set the initial value of $\Delta = 2^{\lfloor \log U \rfloor}$, that is, the largest power of 2 less than or equal to $U = \max\{b(i) : i \in N\}$. In a strongly polynomial algorithm, we cannot take logarithms because we cannot determine $\log U$ in $O(1)$ elementary arithmetic operations. Therefore, in the enhanced capacity scaling algorithm, we set $\Delta = \max\{b(i) : i \in N\}$.
2. The capacity scaling algorithm decreases Δ by a factor of 2 in every scaling phase. In the enhanced capacity scaling algorithm, we also decrease Δ by a factor of 2, but if $\max\{|e(i)| : i \in N\} \leq \Delta/8n$, then we reset $\Delta = \max\{|e(i)| : i \in N\}$. Consequently, the enhanced capacity scaling algorithm generally decreases Δ by a factor of 2, but sometimes by a larger factor when imbalances are too small compared to the current scale factor. Without resetting Δ in this way, the capacity scaling algorithm might perform $O(\log U)$ scaling phases, many of which will not perform any augmentations. The resulting algorithm would contain $O(\log U)$ in its running time and would not be strongly polynomial-time.
3. In the capacity scaling algorithm, each arc flow is an integral multiple of Δ . This property is essential for its correctness because it ensures that each positive residual capacity is a multiple of Δ , and consequently, any augmentation can carry Δ units of flow. In the enhanced capacity scaling algorithm, although the flows on nonabundant arcs are integral multiples of Δ , the flows on the abundant arcs can be arbitrary. Since the flows on abundant arcs are sufficiently large, their arbitrary values do not prohibit sending Δ units of flow on them.
4. The capacity scaling algorithm sends Δ units of flow from a node k with $e(k) \geq \Delta$ to a node l with $e(l) \leq -\Delta$. As a result, the excess nodes do not become

deficit nodes, and vice versa. In the enhanced capacity scaling algorithm, augmentations carry Δ units of flow and are (a) either from a node k with $e(k) > (n - 1)\Delta/n$ to a node l with $e(l) < -\Delta/n$, (b) or from a node k with $e(k) > \Delta/n$ to a node l with $e(l) < -(n - 1)\Delta/n$. Notice that due to these choices, excess nodes might become deficit nodes and deficit nodes might become excess nodes. Although these choices might seem a bit odd when compared to the capacity scaling algorithm, they ensure several nice theoretical properties that we describe in the following discussion.

We establish the correctness of the enhanced capacity scaling algorithm as follows. In the Δ -scaling phase, we refer to a node i as a *large excess node* if $e(i) > (n - 1)\Delta/n$ and as a *medium excess node* if $e(i) > \Delta/n$. (Observe that a large excess node is also a medium excess node.) Similarly, we refer to a node i as a *large deficit node* if $e(i) < -(n - 1)\Delta/n$ and as a *medium deficit node* if $e(i) < -\Delta/n$. In the Δ -scaling phase, each shortest path augmentation either starts at a large excess node k and ends at a medium deficit node l , or starts at a medium excess node k and ends at a large deficit node l . To establish the correctness of the algorithm, we need to show that whenever (1) the network contains a large excess node k , it must also contain a medium deficit node l , or when (2) the network contains a large deficit node l , it must also contain a medium excess node k . We establish this result in the following lemma.

Lemma 10.27. *If the network contains a large excess node k , it must also contain a medium deficit node l . Similarly, if the network contains a large deficit node l , it must also contain a medium excess node k .*

Proof. We prove the first part of the lemma; the proof of the second part is similar. Note that $\sum_{i \in N} e(i) = 0$ because the total excess of the excess nodes equals the total deficit of the deficit nodes. If $e(k) > (n - 1)\Delta/n$ for some excess node k , the total deficit of deficit nodes is also greater than $(n - 1)\Delta/n$. Since the network contains at most $(n - 1)$ deficit nodes, at least one of these nodes, say node l , must have a deficit greater than Δ/n , or equivalently $e(l) < -\Delta/n$. ♦

In the proofs, we use the following lemma several times.

Lemma 10.28. *At the end of the Δ -scaling phase, $|e(i)| \leq (n - 1)\Delta/n$ for each node i . At the beginning of the Δ -scaling phase, $|e(i)| \leq 2(n - 1)\Delta/n$ for each node i .*

Proof. Suppose that during some scaling phase the network contains some large excess node. Then by Lemma 10.27, it also contains some medium deficit node, so the scaling phase would not yet end. Similarly, if the network contains some large deficit node, it would also contain some medium excess node, and the scaling phase would not end. Therefore, at the end of the scaling phase, $|e(i)| \leq (n - 1)\Delta/n$ for each node i .

If at the next scaling phase the algorithm halves the value of Δ , then $|e(i)| \leq 2(n - 1)\Delta/n$ for each node i . On the other hand, if the algorithm sets Δ equal to e_{\max} , then $|e(i)| \leq \Delta$ for each node i . In either case, the lemma is true. ♦

The enhanced capacity scaling algorithm also relies on the fact that in the Δ -scaling phase, we can send Δ units of flow along the shortest path P from node k to node l . To prove this result, we need to show that the residual capacity of every arc in the path P is at least Δ . We establish this property in two parts. First, we show that the flow on each nonabundant arc is a multiple of Δ ; this would imply that residual capacities of nonabundant arcs and their reversals in the residual network are multiples of Δ (because all the arcs in A are uncapacitated). We next show that the flow on each abundant arc is always greater than or equal to $4n\Delta$; therefore, we can send Δ units of flow in either direction. These two results would complete the correctness proof of the enhanced capacity scaling algorithm.

Lemma 10.29. *Throughout the execution of the enhanced capacity scaling algorithm, the solution satisfies the flow and imbalance properties (i.e., Properties 10.25 and 10.26).*

Proof. We prove this lemma by performing induction on the number of flow augmentations and changes in the scale factor Δ . We first consider the flow property. Each augmentation sends Δ units of flow and thus preserves the property. The scale factor Δ changes in one of the two following ways: (1) when we replace Δ by $\Delta' = \Delta/2$, or (2) after replacing $\Delta' = \Delta/2$, we reset $\Delta'' = \max\{|e(i)| : i \in N\}$. In case (1), the flows on the nonabundant arcs continue to be multiples of Δ' . In case (2), $\Delta'' = \max\{|e(i)| : i \in N\} \leq \Delta'/8n$, or $\Delta' \geq 8n\Delta''$. Since each positive arc flow x_{ij} on a nonabundant arc is a multiple of Δ' , $x_{ij} \geq \Delta' \geq 8n\Delta''$. Consequently, each positive flow arc becomes an abundant arc (with respect to the new scale factor) and vacuously satisfies the flow property.

We next establish the imbalance property by performing induction on the number of augmentations and the creation of new abundant arcs. Each augmentation carries flow from a nonroot node to another nonroot node and preserves the property. Moreover, each time the algorithm creates a new abundant arc, it might create a nonroot node i with nonzero imbalance; however, it immediately performs an imbalance-property augmentation to reduce its imbalance to zero. The lemma now follows. ♦

Theorem 10.30. *In the Δ -scaling phase, the algorithm changes the flow on any arc by at most $4n\Delta$ units.*

Proof. The flow on an arc changes through either imbalance-property augmentations or shortest path augmentations. We first consider changes caused by imbalance-property augmentations. At the beginning of the Δ -scaling phase, $|e(i)| \leq 2(n-1)\Delta/n$ for each node i (from Lemma 10.28). Consequently, an imbalance-property augmentation changes the flow on any arc by at most $2(n-1)\Delta/n$. Since the algorithm can perform at most n imbalance-property augmentations at the beginning of a scaling phase, the change in the flow on an arc due to all imbalance-property augmentations is at most $2(n-1)\Delta \leq 2n\Delta$.

Next consider the changes in the flow on an arc caused by shortest path augmentations. At the beginning of the Δ -scaling phase, each root node i satisfies the condition $|e(i)| \leq 2(n-1)\Delta/n$ (by Lemma 10.29). Consider the case when the Δ -scaling phase performs no imbalance-property augmentations. In this case, at most

one shortest path augmentation will begin at a large excess node i , because after this augmentation, the new excess $e'(i)$ satisfies the inequality $e'(i) \leq 2(n-1)\Delta/n - \Delta = (n-2)\Delta/n \leq (n-1)\Delta/n$, and node i is no longer a large excess node. Similarly, at most one shortest path augmentation will end at a large deficit node.

Now suppose that the algorithm does perform some imbalance-property augmentations. In this case the algorithm sends $e(j)$ units of flow from each nonroot node j to the root of its abundant component. The subsequent imbalance-property augmentation from node j to the root node i can increase $|e(i)|$ by at most $2(n-1)\Delta/n$ units, so node i can be the start or end node of at most two additional shortest path augmentations in the Δ -scaling phase. We “charge” these two augmentations to node j , which becomes a nonroot node and remains a nonroot node in the subsequent scaling phases.

To summarize, we have shown that in the Δ -scaling phase, we can charge each root node at most one shortest path augmentation and each nonroot node at most two shortest path augmentations. Each such augmentation changes the flow on any arc by 0 or Δ units. Consequently, the total flow change on any arc due to all shortest path augmentations is at most $2n\Delta$. We have earlier shown the total flow change due to imbalance-property augmentations is at most $2n\Delta$. These results establish the theorem. ♦

The preceding theorem immediately implies the following result.

Lemma 10.31. *If the algorithm designates an arc (i, j) as an abundant arc in the Δ -scaling phase, then in all subsequent Δ' -scaling phases $x_{ij} \geq 4n\Delta'$.*

Proof. We prove this result by performing induction on the number of scaling phases. Since the algorithm designates arc (i, j) as an abundant arc at the beginning of the Δ -scaling phase, the flow on this arc satisfies the condition $x_{ij} \geq 8n\Delta$. The Lemma 10.31 implies that the flow change on any arc in the Δ -scaling phase is at most $4n\Delta$. Therefore, throughout the Δ -scaling phase and, also, at the end of this scaling phase, the arc (i, j) satisfies the condition $x_{ij} \geq 4n\Delta$. In the next scaling phase, the scale factor $\Delta' \leq \Delta/2$; so at the beginning of the Δ' -scaling phase, $x_{ij} \geq 8n\Delta'$. This conclusion establishes the lemma. ♦

We next consider the worst-case complexity of the enhanced capacity scaling algorithm. We show that the algorithm performs $O(n \log n)$ scaling phases, requiring a total of $O(n \log n)$ shortest path augmentations. These proofs rely on the result, stated in Theorem 10.33, that any abundant component whose root node has a medium excess or a medium deficit merges into a larger abundant component within $O(\log n)$ scaling phases. Theorem 10.33, in turn, depends on the following lemma.

Lemma 10.32. *Let S be the set of nodes spanned by an abundant component, and let $e(S) = \sum_{i \in S} e(i)$ and $b(S) = \sum_{i \in S} b(i)$. Then $b(S) - e(S)$ is an integral multiple of Δ .*

Proof. Summing the mass balance constraints (9.1b) of nodes in S , we see that

$$b(S) - e(S) = \sum_{\{(i,j) \in (S, \bar{S})\}} x_{ij} - \sum_{\{(i,j) \in (\bar{S}, S)\}} x_{ij}. \quad (10.12)$$

In this expression, (S, \bar{S}) and (\bar{S}, S) denote the sets of forward and backward arcs in the cut $[S, \bar{S}]$. Since the flow on each arc in the cut is an integral multiple of Δ (by the flow property), $b(S) - e(S)$ is also an integral multiple of Δ . ♦

Theorem 10.33. *Let S be the set of nodes spanned by an abundant component and suppose that at the end of the Δ -scaling phase, $|e(S)| > \Delta/n$. Then within $O(\log n)$ additional scaling phases, the algorithm will merge the abundant component S into a larger abundant component.*

Proof. We first claim that at the end of the Δ -scaling phase, $|b(S)| \geq \Delta/n$. We prove this result by contradiction. Suppose that $|b(S)| < \Delta/n$. Let node i be the root node of the component S . Lemma 10.28 implies that at the end of the Δ -scaling phase, $|e(i)| = |e(S)| \leq (n-1)\Delta/n$. Therefore, $|b(S)| + |e(S)| < \Delta$, which from Lemma 10.32 is possible only if $|b(S)| = |e(S)|$. This condition, however, contradicts the facts that $|e(S)| > \Delta/n$ and $|b(S)| < \Delta/n$. Therefore, $|b(S)| \geq \Delta/n$ whenever $|e(S)| > \Delta/n$. Consequently, at the end of the Δ -scaling phase, $|b(S)| \geq \Delta/n$.

Since the enhanced capacity scaling algorithm decreases Δ by a factor of at least 2 in each scaling phase, within $\log(9n^2m) \leq \log(9n^4) = O(\log n)$ scaling phases, the scale factor will be $\Delta' \leq \Delta/2^{\log(9n^2m)} = \Delta/(9n^2m)$, or $\Delta/n \geq 9nm\Delta'$. Since $|b(S)| \geq \Delta/n$, $|b(S)| \geq 9nm\Delta'$. We consider the situation when $b(S) > 0$. [The analysis of the situation with $b(S) < 0$ is similar.] Since $e(S) \leq \Delta'(n-1)/n \leq \Delta'$ (by Lemma 10.28), the flow across the cut $[S, \bar{S}]$ (i.e., the right-hand side of (10.12)) is at least $9nm\Delta' - \Delta' \geq 8nm\Delta'$. This cut contains at most m arcs; at least one of these arcs, say arc (i, j) , must have a flow at least $8n\Delta'$. Thus the algorithm will designate the arc (i, j) as an abundant arc and merge the component S into a larger abundant component. ♦

We are now ready to complete the proof of the main result of this section.

Theorem 10.34. *The enhanced capacity scaling algorithm solves the uncapacitated minimum cost flow problem within $O(n \log n)$ scaling phases and performs a total of $O(n \log n)$ shortest path augmentations. If $S(n, m)$ is the time required to solve a shortest path problem with nonnegative arc lengths, the running time of the enhanced capacity scaling algorithm is $O(n \log n S(n, m))$.*

Proof. We first show that the algorithm performs $O(n \log n)$ scaling phases. Consider a scaling phase with scale factor equal to Δ . At the end of this scaling phase, we will encounter one of the following two outcomes:

Case 1. For some node i , $|e(i)| > \Delta/16n$. Let node i be the root node of an abundant component S . Clearly, within four scaling phases, either the component S merges into a larger component or $|e(i)| > \Delta/n$. In the latter case, Theorem 10.33 implies that within $O(\log n)$ scaling phases, the component S merges into a larger component.

Case 2. For every node i , $|e(i)| \leq \Delta/16n$. At the beginning of the next scaling phase, the new scale factor $\Delta' = \Delta/2$, so $|e(i)| \leq \Delta'/8n$ for each node i . We then reset $\Delta' = \max\{|e(i)| : i \in N\}$. As a result, for some node i , $|e(i)| =$

$\Delta' > \Delta'/16n$ and, as in Case 1, within $O(\log n)$ scaling phases, the abundant component containing node i merges into a larger component.

This discussion shows that within $O(\log n)$ scaling phases the algorithm performs one merge operation. Since each merge operation decreases the number of abundant components by one, the algorithm can perform at most n merge operations. Consequently, the number of scaling phases is bounded by $O(n \log n)$. The algorithmic description of the enhanced capacity scaling algorithm implies that the algorithm requires $O(m)$ time per scaling phase plus the time required for the augmentations.

We now obtain a bound on the number of augmentations and the time that they require. The algorithm performs at most n imbalance-property augmentations; it can easily execute each augmentation in $O(m)$ time; thus these augmentations are not a bottleneck step in the algorithm. Next consider the shortest path augmentations. Recall from the proof of Theorem 10.30 that in a scaling phase, we can charge each shortest path augmentation to a root node (which is a large excess or a large-deficit node) or to a nonroot node. Since we can charge each nonroot at most two augmentations over the entire execution of the algorithm, we charge at most $2n$ augmentations to nonroots. Moreover, when we charge an augmentation to a root node i , this node satisfies the condition $|e(i)| \geq (n-1)\Delta/n$. Theorem 10.33 implies that we will charge at most one augmentation to node i in the following $O(\log n)$ scaling phases before the algorithm performs a merge operation and the component containing node i merges into a larger component. Since the algorithm encounters at most $2n$ different abundant components (n to begin with and n due to merge operations), the total number of shortest path augmentations we can charge to root nodes is at most $O(n \log n)$. Since each shortest path augmentation requires the solution of a shortest path problem with nonnegative arc lengths and requires $S(n, m)$ time, all the shortest path augmentations require a total of $O(n \log n S(n, m))$ time. This time dominates the time taken by all other operations performed by the algorithm. Therefore, we have established the assertion of the theorem. ♦

To solve the capacitated minimum cost flow problem, we transform it to the uncapacitated version using the transformation described in Section 2.4. The resulting uncapacitated network has $n' = n + m$ nodes and $m' = 2m$ arcs. The enhanced capacity scaling algorithm will solve the minimum cost flow problem in the transformed network in $O(n' \log n') = O(m \log m) = O(m \log n^2) = O(m \log n)$ scaling phases and will solve a total of $O(n' \log n') = O(m \log n)$ shortest path problems. Each shortest path problem in the uncapacitated network requires $S(n', m')$ time, but using the ideas described in Exercise 4.53 we can improve this time bound to $S(n, m)$. Therefore, the enhanced capacity scaling algorithm can solve the capacitated minimum cost flow problem in $O(m \log n S(n, m))$ time. We state this important result as a theorem.

Theorem 10.35. *The enhanced capacity scaling algorithm solves a capacitated minimum cost flow problem in $O(m \log n S(n, m))$ time.* ♦

10.8 SUMMARY

In this chapter we continued our study of the minimum cost flow problem by developing several polynomial-time algorithms. The scaling technique is a central theme in almost all the algorithms we have discussed. The algorithms discussed use capacity scaling, cost scaling, or both, or use scaling concepts in their proofs. We discussed six polynomial-time algorithms: (1) the capacity scaling algorithm, (2) the cost scaling algorithm, (3) the double scaling algorithm, (4) the minimum mean cycle-canceling algorithm, (5) the repeated capacity scaling algorithm, and (6) the enhanced capacity scaling algorithm. The first three of these algorithms are weakly polynomial; the other three are strongly polynomial. Figure 10.13 specifies the running times of these algorithms.

The capacity scaling algorithm is possibly the simplest of all the polynomial-time algorithms we have discussed. This algorithm is an improved version of the successive shortest path algorithm discussed in Section 9.7; by augmenting flows along paths with sufficiently large residual capacities, this algorithm is able to decrease the number of augmentations from $O(nU)$ to $O(m \log U)$.

Whereas the capacity scaling algorithm scales the capacities, the cost scaling algorithm scales costs. The algorithm maintains ϵ -optimal flows for decreasing values of ϵ and repeatedly executes an improve-approximation procedure that converts an ϵ -optimal flow into an $\epsilon/2$ -optimal flow. The computations performed by the improve-approximation procedure are similar to those performed by the preflow-push algorithm for the maximum flow problem. The double scaling algorithm is the same as the cost scaling algorithm except that it uses a different version of the improve-approximation procedure. The improve-approximation procedure in the cost scaling algorithm performs push/relabel steps; in the double scaling algorithm, this procedure augments flow along paths of sufficiently large residual capacity. Justifying its name, within a cost scaling phase, the double scaling algorithm performs a number of capacity scaling phases.

The minimum mean cycle-canceling algorithm for the minimum cost flow problem is different from all the other algorithms discussed in this chapter. The algorithm is startlingly simple to describe and does not make explicit use of the scaling technique; the proof of the algorithm, however, uses arguments from scaling techniques.

Algorithm	Running time
Capacity scaling algorithm	$O((m \log U)(m + n \log n))$
Cost scaling algorithm	$O(n^3 \log(nC))$
Double scaling algorithm	$O(nm \log U \log(nC))$
Minimum mean cycle-canceling algorithm	$O(n^2 m^3 \log n)$
Repeated capacity scaling algorithm	$O((m^2 \log n)(m + n \log n))$
Enhanced capacity scaling algorithm	$O((m \log n)(m + n \log n))$

Figure 10.13 Running times of polynomial-time minimum cost flow algorithms.

This algorithm is a special implementation of the cycle canceling algorithm that we described in Section 9.6; it always augments flow along a minimum mean (negative) cycle in the residual network. To establish that this algorithm is strongly polynomial, we show that (1) when the reduced cost of an arc is *sufficiently large*, the flow on the arc becomes “fixed” (i.e., does not change any more); and (2) within $O(nm \log n)$ iterations, at least one additional arc has a sufficiently large reduced cost so that its value becomes fixed.

If we adopt a similar idea in the capacity scaling algorithm, it also becomes strongly polynomial. We showed that whenever the flow on an arc (i, j) is sufficiently large, we can fix the potentials of nodes i and j with respect to each other. The repeated capacity scaling algorithm applies the capacity scaling algorithm and within $O(\log n)$ scaling phases, it identifies an arc (i, j) with a sufficiently large flow. The algorithm then merges the nodes i and j into a single node and starts from scratch again on the modified minimum cost flow problem. The enhanced capacity scaling algorithm, described next, dramatically improves on the repeated capacity scaling algorithm by observing that whenever we contract an arc, we need not start all over again, but can continue the computations and still contract an additional arc within every $O(\log n)$ scaling phases and use only $O(m \log n)$ augmentations in total. This algorithm does not perform contractions explicitly, but does so implicitly by maintaining zero excesses at the contracted nodes (i.e., nonroot nodes).

REFERENCE NOTES

The following account of polynomial-time minimum cost flow algorithms is fairly brief. The surveys by Ahuja, Magnanti, and Orlin [1989, 1991] and by Goldberg, Tardos, and Tarjan [1989] provide more details concerning the development of this field.

Most of the available (combinatorial) polynomial-time algorithms for the minimum cost flow problems use scaling techniques. Edmonds and Karp [1972] introduced the scaling approach and obtained the first weakly polynomial-time algorithm for the minimum cost flow problem. This algorithm used the capacity scaling technique. The algorithm we presented in Section 10.2, which is a variant of Edmonds and Karp’s algorithm, is due to Orlin [1988]. From 1972 to 1984, there was little research on scaling techniques. Since 1985, research employing scaling techniques has been extensive. Researchers now recognize that scaling techniques have great theoretical value as well as potential practical significance. Scaling techniques now yield many of the best (in the worst-case sense) available minimum cost flow algorithms.

Röck [1980] and, independently, Bland and Jensen [1985] suggested a cost scaling technique for the minimum cost flow problem. This approach solves the minimum cost flow problem as a sequence of $O(n \log C)$ maximum flow problems. Goldberg and Tarjan [1987] improved on the running time of Röck’s algorithm and solved the minimum cost flow problem by solving “almost” $O(\log(nC))$ maximum flow problems. This approach is based on the concept of ϵ -optimality, which is, independently, due to Bertsekas [1979] and Tardos [1985]. We describe this approach in Section 10.3. Goldberg and Tarjan [1987] have developed several improved implementations of this approach, including the wave implementation presented in

Section 10.3. Their best implementation, which runs in $O(nm \log(n^2/m) \log(nC))$ time, uses Fibonacci heaps and finger search trees. Bertsekas and Eckstein [1988], independently, discovered the wave implementation.

Ahuja, Goldberg, Orlin, and Tarjan [1992] developed the double scaling algorithm described in Section 10.4, which combines capacity and cost scaling. This paper also describes several improved implementations, the best of which runs in $O(nm \log \log U \log(nC))$ time and uses the Fibonacci heap data structure.

When Edmonds and Karp [1972] suggested the first (weakly) polynomial-time algorithm for the minimum cost flow problem, they posed the development of a strongly polynomial-time algorithm as an open challenging problem. Tardos [1985] first settled this problem. Subsequently, Orlin [1984], Fujishige [1986], Galil and Tardos [1986], Goldberg and Tarjan [1987, 1988], Orlin [1988], and Ervolina and McCormick [1990b] developed other strongly polynomial-time algorithms. Currently, the best strongly polynomial-time algorithm is due to Orlin [1988]; it runs in $O((m \log n)(m + n \log n))$ time.

Most of the strongly polynomial-time minimum cost flow algorithm use the ideas of “fixing arc flows” or “fixing node potentials.” Tardos [1985] was the first investigator to propose the use of either of these ideas (her algorithm fixes arc flows). The minimum mean cycle-canceling algorithm that we presented in Section 10.5 fixes arc flows; it is due to Goldberg and Tarjan [1988]. Goldberg and Tarjan [1988] also presented several variants of the minimum mean cycle-canceling algorithm with improved worst-case complexity. Orlin [1984] and Fujishige [1986] independently developed the idea of fixing node potentials, which is the “dual” of fixing arc flows. Using this idea, Goldberg, Tardos, and Tarjan [1989] obtained the repeated capacity scaling algorithm that we examined in Section 10.6. The enhanced capacity scaling algorithm, which is due to Orlin [1988], achieves the best strongly polynomial-time for solving the minimum cost flow problem. However, our presentation of the enhanced capacity scaling algorithm in Section 10.7 is based on Plotkin and Tardos’ [1990] simplification of Orlin’s original algorithm.

Some additional polynomial-time minimum cost flow algorithms include (1) a triple scaling algorithm due to Gabow and Tarjan [1989a], (2) a special implementation of the cycle canceling algorithm developed by Barahona and Tardos [1989], and (3) (its dual approach) a cut canceling algorithm proposed by Ervolina and McCormick [1990a].

Interior point linear programming algorithms are another source of polynomial-time algorithms for the minimum cost flow problem. Among these, the fastest available algorithm, due to Vaidya [1989], solves the minimum cost flow problem in $O(n^{2.5} \sqrt{m} K)$ time, with $K = \log n + \log C + \log U$.

Currently, the best available time bound for the minimum cost flow problem is $O(\min\{nm \log(n^2/m) \log(nC), nm (\log \log U) \log(nC), (m \log n)(m + n \log n)\})$; the three bounds in this expression are, respectively, due to Goldberg and Tarjan [1987], Ahuja, Goldberg, Orlin, and Tarjan [1992], and Orlin [1988].

EXERCISES

- 10.1. Suppose that we want to solve the minimum cost flow problem shown in Figure 10.14(a) by the capacity scaling algorithm. Show the computations for two scaling phases. You may identify the shortest path distances by inspection.

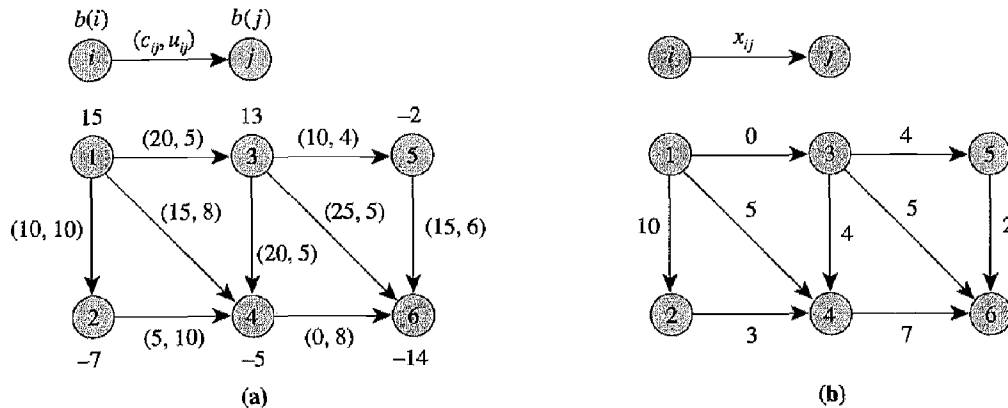


Figure 10.14 Examples for Exercises 10.1 and 10.4.

- 10.2. In every iteration of the capacity scaling algorithm, we augment flow along a shortest path from a node k with $e(k) \geq \Delta$ to a node l with $e(l) \leq -\Delta$. Suppose that we modify the algorithm as follows: We let node l be any deficit node; that is, we do not necessarily assume that $e(l) \leq -\Delta$. Will this modification affect the worst-case complexity of the capacity scaling algorithm?
- 10.3. Prove or disprove the following statements.
- During the Δ -scaling phase of the capacity scaling algorithm, $|e(i)| \leq 2\Delta$ for each node $i \in N$.
 - While solving a specific instance of the minimum cost flow problem, the capacity scaling algorithm might perform more augmentations than the successive shortest path algorithm.
- 10.4. Consider the minimum cost flow problem given in Figure 10.14(a) and the feasible flow x shown in Figure 10.14(b). Starting with $\epsilon = 0$, apply two phases of the cost scaling algorithm.
- 10.5. Show that if the cost-scaling algorithm finds that arc (i, j) is inadmissible at some stage, this arc remains inadmissible until the algorithm relabels node i .
- 10.6. Let x and x' be two distinct (feasible) flows in a network. The flow decomposition theorem implies that we can always express x' as x plus the flow along at most m directed cycles W_1, W_2, \dots, W_p in $G(x)$. For every $1 \leq i \leq p$, let W'_i denote the directed cycle obtained by reversing each arc in W_i . Show that we can express x as x' plus the flow along the cycles W'_1, W'_2, \dots, W'_p .
- 10.7. For the cost scaling algorithm, we showed that whenever $\epsilon < 1/n$, any ϵ -optimal flow is 0-optimal. Show that if we multiply all arc costs by $n + 1$, then any flow that is ϵ -optimal flow for the modified problem when $\epsilon \leq 1$ is 0-optimal for the original problem.
- 10.8. In the cost scaling algorithm, during a relabel operation we increase node potentials by $\epsilon/2$ units. Show that we can increase node potentials by as much as $\epsilon/2 + \min\{c_{ij}^T : (i, j) \in G(x) \text{ and } r_{ij} > 0\}$ and still maintain $\epsilon/2$ -optimality of the pseudoflow.
- 10.9. Let x' be a feasible flow of the minimum cost flow problem and let x be a pseudoflow. Show that in the pseudoflow x , for every node v with an excess, there exists a node w with a deficit and a sequence of nodes $v = v_0, v_1, v_2, \dots, v_l = w$ that satisfies the property that the path $P = v_0 - v_1 - v_2 - \dots - v_l$ is a directed path in $G(x)$ and its reversal $\bar{P} = v_l - v_{l-1} - \dots - v_0$ is a directed path in $G(x')$. (Hint: This exercise is similar to Exercise 10.6.)
- 10.10. In this exercise we study the nonscaled version of the cost scaling algorithm.
- Modify the algorithm described in Section 10.3 so that it starts with a 0-optimal

pseudoflow, maintains an $1/(n + 1)$ -optimal pseudoflow at every step, and terminates with an $1/(n + 1)$ -optimal flow.

- (b) Determine the number of relabel operations, the number of saturating and non-saturating pushes, and the running time of the algorithm. Compare these numbers with those of the cost scaling algorithm.
- 10.11. In the wave implementation of the cost scaling algorithm described in Section 10.3, we scaled costs by a factor of 2. Suppose, instead, that we scaled costs by a factor of $k \geq 2$. In that case we start with $\Delta = k^{\lceil \log C \rceil}$ and decrease ϵ by a factor of k between two consecutive scaling phases. Outline the changes required in the algorithm and determine the number of scaling phases, relabel operations, and saturating and non-saturating pushes within a scaling phase. For what value of k is the running time minimum?
- 10.12. **Generalized cost scaling algorithm** (Goldberg and Tarjan [1987]). As we noted in the text, by using some of the ideas of the minimum mean cycle-canceling algorithm (described in Section 10.5), we can devise a strongly polynomial-time version of the cost scaling algorithm that we described in Section 10.3. The modified algorithm, which we call the *generalized cost scaling algorithm*, is the same as the cost scaling algorithm except that it performs the following additional step after it has called the procedure *improve-approximation*, but before resetting $\epsilon : = \epsilon/2$ (see Figure 10.3).

Additional step: Solve a minimum mean cycle problem to determine the minimum mean cycle cost $\mu(x)$, set $\epsilon = -\mu(x)$, and then determine a set of potential π so that the flow x is ϵ -optimal with respect to π (as described in the proof of Lemma 10.12).

Show that the generalized cost scaling fixes a distinct arc after $O(\log n)$ scaling phases. What is the resulting running time of the algorithm?

- 10.13. In the double scaling algorithm described in Section 10.4, we scaled costs by a factor of 2. Suppose that as described in Exercise 10.2, we scale costs by a factor of k instead of 2. Show that within a cost scaling phase, the algorithm performs $O(knm)$ retreat steps. How many advance steps does the algorithm perform within a scaling phase? How many scaling phases does it require? For what value of k does the algorithm run in the least time? What is the time bound for this value of k ?
- 10.14. An arc (i, j) in the network $G = (N, A)$ is *critical* if increasing c_{ij} causes the cost of the optimal flow to increase and decreasing c_{ij} causes the cost of the optimal flow to decrease. Does a network always contain a critical arc? Show that we can identify all critical arcs by solving $O(m)$ maximum flow problems. (*Hint:* Use the fact that an arc is critical if it carries a positive flow in every optimal flow.)
- 10.15. In some minimum cost flow problem, each arc capacity and each supply/demand is a multiple of α and lies in the range $[0, \alpha K]$ for some constant K . Will the algorithms discussed in this chapter run any faster when applied to minimum cost flow problems with this special structure?
- 10.16. Suppose that in some minimum cost flow problem, each arc cost is a multiple of α and lies in the range $[0, \alpha K]$ for some constant K . Will this special structure permit us to solve the minimum cost flow problem any faster by the cost scaling and double scaling algorithms?
- 10.17. **Minimum cost flows in unit capacity networks.** A network is a *unit capacity network* if each arc has a capacity of 1.
 - (a) What is the running time of the capacity scaling algorithm for unit capacity networks?
 - (b) What is the running time of the cost scaling algorithm for unit capacity networks? (*Hint:* Will the algorithm make any nonsaturating pushes?)
- 10.18. **Minimum cost flows in bipartite networks.** Let $G = (N_1 \cup N_2, A)$ be a bipartite network. Let $n_1 = |N_1| \leq |N_2| = n_2$.
 - (a) Show that when applied to a bipartite network, the cost scaling algorithm relabels any node $O(n_1)$ times during a scaling phase.

- (b) Develop an implementation of the generic cost scaling algorithm that runs in $O(n_1^2 m \log(nC))$ time for bipartite networks. (*Hint*: Generalize the bipartite preflow-push algorithm for the maximum flow problem discussed in Section 8.3.)
- 10.19. What is the running time of the double scaling algorithm for bipartite networks $G = (N_1 \cup N_2, A)$, assuming that $n_1 = |N_1| \leq |N_2| = n_2$?
- 10.20. Two minimum cost flow problems P' and P'' are *capacity adjacent* if P'' differs from P' only in one arc capacity and by 1 unit. Given an optimal solution of P' , describe an efficient method for solving P'' . (*Hint*: Reoptimize by solving a shortest path problem.)
- 10.21. Two minimum cost flow problems P' and P'' are *cost adjacent* if P'' differs from P' only in one arc cost, and by 1 unit. Given an optimal solution of P' , describe an efficient method for solving P'' . (*Hint*: Reoptimize by solving a maximum flow problem.)
- 10.22. **Bit scaling of capacities** (Röck [1980]). In this capacity scaling algorithm, we consider binary representations of the arc capacities (as described in Section 3.3) and define problem P^k to be the minimum cost flow problem with each arc capacity equal to the k leading bits of the actual capacity. Given an optimal solution of P^k , how would you obtain an optimal solution of P^{k+1} by solving at most m capacity adjacent problems (as defined in Exercise 10.20). Write a pseudocode for the minimum cost flow problem assuming the availability of a subroutine for solving capacity adjacent problems (i.e., solving one from the solution to the other). What is the running time of your algorithm?
- 10.23. **Bit scaling of costs** (Röck [1980]). In this cost scaling algorithm, we consider binary representations of the arc costs and define problem P^k to be the minimum cost flow problem with each arc cost equal to the k leading bits of the actual cost. Given an optimal solution of P^k , how would you obtain an optimal solution of P^{k+1} by solving at most m cost adjacent problems (as defined in Exercise 10.21)? Write a pseudocode for the minimum cost flow problem assuming the availability of a subroutine for solving cost adjacent problems (i.e., solving one from the solution to the other). What is the running time of your algorithm?
- 10.24. Suppose that we define the contraction of an arc as in Section 10.5. Let G^c denote the network of $G = (N, A)$ we obtain when we contract the endpoints of an arc $(k, l) \in A$ into a single node p . In addition, let $G' = (N, A - \{(k, l)\})$. Show that if $\alpha(G)$ denotes the number of (distinct) spanning trees of G , then $\alpha(G) = \alpha(G^c) + \alpha(G')$.
- 10.25. **Constrained maximum flow problem.** In the constrained maximum flow problem, we wish to maximize the flow from the source node s to the sink node t subject to an additional linear constraint. Consider the following linear programming formulation of this problem:

Maximize v

subject to

$$\sum_{\{j: (i,j) \in A\}} x_{ij} - \sum_{\{j: (j,i) \in A\}} x_{ji} = \begin{cases} v & \text{for } i = s \\ 0 & \text{for all } i \in N - \{s, t\} \\ -v & \text{for } i = t, \end{cases}$$

$$0 \leq x_{ij} \leq u_{ij},$$

$$\sum_{(i,j) \in A} c_{ij} x_{ij} \leq D.$$

- (a) Let v^* be any integer and let x^* be an optimal solution of a minimum cost flow problem with the objective function $\sum_{(i,j) \in A} c_{ij} x_{ij}$ and with the supply/demand data $b(s) = v^*$, $b(t) = -v^*$, and $b(i) = 0$ for all other nodes. Let $z^* = \sum_{(i,j) \in A} c_{ij} x_{ij}^*$. Show that x^* solves the constrained maximum flow problem when $D = z^*$. Assume that $c_{ij} \geq 0$ for each arc $(i, j) \in A$.
- (b) Assume that all of the data in the constrained maximum flow problem are integer. Use the result in part (a) to develop an algorithm for the constrained maximum

flow problem that uses a minimum cost flow algorithm as a subroutine. What is the running time of your algorithm? (*Hint*: Perform binary search on v .)

- 10.26.** In the enhanced capacity scaling algorithm, suppose we maintain an index with each arc that stores whether the arc is an abundant or a nonabundant arc. Suppose further that at some stage the algorithm adds an arc (i, j) to the abundant subgraph. Show how you would perform each of the following operations in $O(m)$ time: (i) identifying the root nodes, i_r and j_r , of the abundant components containing the nodes i and j ; (ii) determining whether the nodes i and j belong to the same abundant component; and (iii) identifying a path from node i to j , or vice versa. Using these operations, explain how you would perform a merge operation and the subsequent imbalance-property augmentation in $O(m)$ time. (*Hint*: Observe that each abundant arc can be traversed in either direction because it has sufficient residual capacity in both the directions. Then use the search algorithm described in Section 3.4.)

11

MINIMUM COST FLOWS: NETWORK SIMPLEX ALGORITHMS

*. . . seek, and ye shall find.
—The Book of Matthew*

Chapter Outline

11.1	Introduction
11.2	Cycle Free and Spanning Tree Solutions
11.3	Maintaining a Spanning Tree Structure
11.4	Computing Node Potentials and Flows
11.5	Network Simplex Algorithm
11.6	Strongly Feasible Spanning Trees
11.7	Network Simplex Algorithm for the Shortest Path Problem
11.8	Network Simplex Algorithm for the Maximum Flow Problem
11.9	Related Network Simplex Algorithms
11.10	Sensitivity Analysis
11.11	Relationship to Simplex Method
11.12	Unimodularity Property
11.13	Summary

11.1 INTRODUCTION

The simplex method for solving linear programming problems is perhaps the most powerful algorithm ever devised for solving constrained optimization problems. Indeed, many members of the academic community view the simplex method as not only one of the principal computational engines of applied mathematics, computer science, and operations research, but also as one of the landmark contributions to computational mathematics of this century. The algorithm has achieved this lofty status because of the pervasiveness of its applications throughout many problem domains, because of its extraordinary efficiency, and because it permits us to not only solve problems numerically, but also to gain considerable practical and theoretical insight through the use of sensitivity analysis and duality theory.

Since minimum cost flow problems define a special class of linear programs, we might expect the simplex method to be an attractive solution procedure for solving many of the problems that we consider in this text. Then again, because network flow problems have considerable special structure, we might also ask whether the simplex method could possibly compete with other “combinatorial” methods, such as the many variants of the successive shortest path algorithm, that exploit the underlying network structure. The general simplex method, when implemented in

a way that does not exploit underlying network structure, is not a competitive solution procedure for solving minimum cost flow problems. Fortunately, however, if we interpret the core concepts of the simplex method appropriately as network operations, we can adapt and streamline the method to exploit the network structure of the minimum cost flow problem, producing an algorithm that is very efficient. Our purpose in this chapter is to develop this network-based implementation of the simplex method and show how to apply it to the minimum cost flow problem, the shortest path problem, and the maximum flow problem.

We could adopt several different approaches for presenting this material, and each has its own merits. For example, we could start by describing the simplex method for general linear programming problems and then show how to adapt the method for minimum cost flow problems. This approach has the advantage of placing our development in the broader context of more general linear programs. Alternatively, we could develop the network simplex method directly in the context of network flow problems as a particular type of augmenting cycle algorithm. This approach has the advantage of not requiring any background in linear programming and of building more directly on the concepts that we have developed already. We discuss both points of view. Throughout most of this chapter we adopt the network approach and derive the network simplex algorithm from the first principles, avoiding the use of linear programming in any direct way. Later, in Section 11.11, we show that the network simplex algorithm is an adaptation of the simplex method.

The central concept underlying the network simplex algorithm is the notion of spanning tree solutions, which are solutions that we obtain by fixing the flow of every arc not in a spanning tree either at value zero or at the arc's flow capacity. As we show in this chapter, we can then solve uniquely for the flow on all the arcs in the spanning tree. We also show that the minimum cost flow problem always has at least one optimal spanning tree solution and that it is possible to find an optimal spanning tree solution by "moving" from one such solution to another, at each step introducing one new nontree arc into the spanning tree in place of one tree arc. This method is known as the network simplex algorithm because spanning trees correspond to the so-called basic feasible solutions of linear programming, and the movement from one spanning tree solution to another corresponds to a so-called pivot operation of the general simplex method. In Section 11.11 we make these connections.

In the first three sections of this chapter we examine several fundamental ideas that either motivate the network simplex method or underlie its development. In Section 11.2 we show that the minimum cost flow problem always has at least one spanning tree solution. We also show how the network optimality conditions that we have used repeatedly in previous chapters specialize when applied to any spanning tree solution. In keeping with our practice in previous chapters, we use these conditions to assess whether a candidate solution is optimal and, if not, how to modify it to construct a better spanning tree solution.

To implement the network simplex algorithm efficiently we need to develop a method for representing spanning trees conveniently in a computer so that we can perform the basic operations of the algorithm efficiently and so that we can efficiently manipulate the computer representation of a spanning tree structure from step to step. We describe one such approach in Section 11.3.

In Section 11.4 we show how to compute the arc flows corresponding to any spanning tree and associated node potentials so that we can assess whether the particular spanning tree is optimal. These operations are essential to the network simplex algorithm, and since we need to make these computations repeatedly as we move from one spanning tree to another, we need to be able to implement these operations very efficiently. Section 11.5 brings all these pieces together and describes the network simplex algorithm.

In the context of applying the network simplex algorithm and establishing that the algorithm properly solves any given minimum cost flow problem, we need to address a technical issue known as degeneracy (which occurs when one of the arcs in a spanning tree, like the nontree arcs, has a flow value equal to zero or the arc's flow capacity). In Section 11.6 we describe a very appealing and simple way to modify the basic network simplex algorithm so that it overcomes the difficulties associated with degeneracy.

Since the shortest path and maximum flow problems are special cases of the minimum cost flow problem, the network simplex algorithm applies to these problems as well. In Sections 11.7 and 11.8 we describe these specialized implementations. When applied to the shortest path problem, the network simplex algorithm closely resembles the label-correcting algorithms that we discussed in Chapter 5. When applied to the maximum flow problem, the algorithm is essentially an augmenting path algorithm.

The network simplex algorithm maintains a feasible solution at each step; by moving from one spanning tree solution to another, it eventually finds a spanning tree solution that satisfies the network optimality conditions. Are there other spanning tree algorithms that iteratively move from one infeasible spanning tree solution to another and yet eventually find an optimal solution? In Section 11.9 we describe two such algorithms: a *parametric network simplex algorithm* that satisfies all of the optimality conditions except the mass balance constraints at two nodes, and a *dual network simplex algorithm* that satisfies the mass balance constraints at all the nodes but might violate the arc flow bounds. These algorithms are important because they provide alternative solution strategies for solving minimum cost flow problems; they also illustrate the versatility of spanning tree manipulation algorithms for solving network flow problems.

We next consider a key feature of the optimal spanning tree solutions generated by the network simplex algorithm. In Section 11.10 we show that it is easy to use these solutions to conduct sensitivity analysis: that is, to determine a new solution if we change any cost coefficient or change the capacity of any arc. This type of information is invaluable in practice because problem data are often only approximate and/or because we would like to understand how robust a solution is to changes in the underlying data.

To conclude this chapter we delineate connections between the network simplex algorithm and more general concepts in linear and integer programming. In Section 11.11 we show that the network simplex algorithm is a special case of the simplex method for general linear programs, although streamlined to exploit the special structure of network flow problems. In particular, we show that spanning trees for the network flow problem correspond in a one-to-one fashion with bases of the linear programming formulation of the problem. We also show that each of

the essential steps of the network simplex algorithm, for example, determining node potentials or moving from one spanning tree to another, are specializations of the usual steps of the simplex method for solving linear programs.

As we have noted in Section 9.6, network flow problems satisfy one very remarkable property: They have optimal integral flows whenever the underlying data are integral. In Section 11.12 we show that this integrality result is a special case of a more general result in linear and integer programming. We define a set of linear programming problems with special constraint matrices, known as *unimodular matrices*, and show that these linear programs also satisfy the integrality property. That is, when solved as linear programs with integral data, problems with these specialized constraint matrices always have integer solutions. Since node–arc incidence matrices satisfy the unimodularity property, this integrality property for linear programming is a strict generalization of the integrality property of network flows. This result provides us with another way to view the integrality property of network flows; it is also suggestive of more general results in integer programming and shows how network flow results have stimulated more general investigations in combinatorial optimization and integer programming.

11.2 CYCLE FREE AND SPANNING TREE SOLUTIONS

Much of our development in previous chapters has relied on a simple but powerful algorithmic idea: To generate an improving sequence of solutions to the minimum cost flow problem, we iteratively augment flows along a series of negative cycles and shortest paths. As one of these variants, the network simplex algorithm uses a particular strategy for generating negative cycles. In this section, as a prelude to our discussion of the method, we introduce some basic background material. We begin by examining two important concepts known as *cycle free solutions* and *spanning tree solutions*.

For any feasible solution, x , we say that an arc (i, j) is a *free arc* if $0 < x_{ij} < u_{ij}$ and is a *restricted arc* if $x_{ij} = 0$ or $x_{ij} = u_{ij}$. Note that we can both increase and decrease flow on a free arc while honoring the bounds on arc flows. However, in a restricted arc (i, j) at its lower bound (i.e., $x_{ij} = 0$) we can only increase the flow. Similarly, for flow on a restricted arc (i, j) at its upper bound (i.e., $x_{ij} = u_{ij}$) we can only decrease the flow. We refer to a solution x as a *cycle free solution* if the network contains no cycle composed only of free arcs. Note that in a cycle free solution, we can augment flow on any augmenting cycle in only a single direction since some arc in any cycle will restrict us from either increasing or decreasing that arc's flow. We also refer to a feasible solution x and an associated spanning tree of the network as a *spanning tree solution* if every nontree arc is a restricted arc. Notice that in a spanning tree solution, the tree arcs can be free or restricted. Frequently, when we refer to a spanning tree solution, we do not explicitly identify the associated tree; rather, it will be understood from the context of our discussion.

In this section we establish a fundamental result of network flows: minimum cost flow problems always have optimal cycle free and spanning tree solutions. The network simplex algorithm will exploit this result by restricting its search for an optimal solution to only spanning tree solutions. To illustrate the argument used to prove these results, we use the network example shown in Figure 11.1.

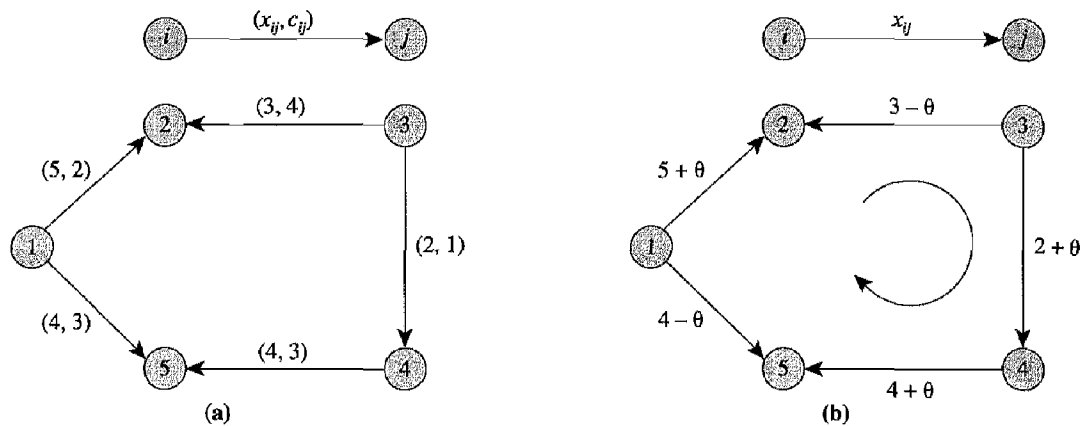


Figure 11.1 Improving flow around a cycle: (a) feasible solution; (b) solution after augmenting θ amount of flow along a cycle.

For the time being let us assume that all arcs are uncapacitated [i.e., $u_{ij} = \infty$ for each $(i, j) \in A$]. The network shown in Figure 11.1 contains positive flow around a cycle. We define the orientation of the cycle as the same as that of arc $(4, 5)$. Let us augment θ units of flow along the cycle in the direction of its orientation. As shown in Figure 11.1, this augmentation increases the flow on arcs along the orientation of the cycle (i.e., forward arcs) by θ units and decreases the flow on arcs opposite to the orientation of the cycle (i.e., backward arcs) by θ units. Also note that the per unit incremental cost for this flow change is the sum of the costs of forward arcs minus the sum of the costs of backward arcs in the cycle, that is,

$$\text{per unit change in cost } \Delta = 2 + 1 + 3 - 4 - 3 = -1.$$

Since augmenting flow in the cycle decreases the cost, we set θ as large as possible while preserving nonnegativity of all arc flows. Therefore, we must satisfy the inequalities $3 - \theta \geq 0$ and $4 - \theta \geq 0$, and hence we set $\theta = 3$. Note that in the new solution (at $\theta = 3$), some arc in the cycle has a flow at value zero, and moreover, the objective function value of this solution is strictly less than the value of the initial solution.

In our example, if we change c_{12} from 2 to 5, the per unit cost of the cycle is $\Delta = 2$. Consequently, to improve the cost by the greatest amount, we would decrease θ as much as possible (i.e., satisfy the restrictions $5 + \theta \geq 0$, $2 + \theta \geq 0$, and $4 + \theta \geq 0$, or $\theta \geq -2$) and again find a lower cost solution with the flow on at least one arc in the cycle at value zero. We can restate this observation in another way: To preserve nonnegativity of all the arc flows, we must select θ in the interval $-2 \leq \theta \leq 3$. Since the objective function depends linearly on θ , we optimize it by selecting $\theta = 3$ or $\theta = -2$, at which point one arc in the cycle has a flow value of zero.

We can extend this observation in several ways:

1. If the per unit cycle cost $\Delta = 0$, we are indifferent to all solutions in the interval $-2 \leq \theta \leq 3$ and therefore can again choose a solution as good as the original one, but with the flow of at least one arc in the cycle at value zero.
2. If we impose upper bounds on the flow (e.g., such as 6 units on all arcs), the

range of flow that preserves feasibility (i.e., the mass balance constraints, lower and upper bounds on flows) is again an interval, in this case $-2 \leq \theta \leq 1$, and we can find a solution as good as the original one by choosing $\theta = -2$ or $\theta = 1$. At these values of θ , the solution is cycle free; that is, some arc on the cycle has a flow either at value zero (at the lower bound) or at its upper bound.

In general, our prior observations apply to any cycle in a network. Therefore, given any initial flow we can apply our previous argument repeatedly, one cycle at a time, and establish the following fundamental result.

Theorem 11.1 (Cycle Free Property). *If the objective function of a minimum cost flow problem is bounded from below over the feasible region, the problem always has an optimal cycle free solution.* ♦

It is easy to convert a cycle free solution into a spanning tree solution. Our results in Section 2.2 show that the free arcs in a cycle free solution define a forest (i.e., a collection of node-disjoint trees). If this forest is a spanning tree, the cycle free solution is already a spanning tree solution. However, if this forest is not a spanning tree, we can add some restricted arcs and produce a spanning tree.

Figure 11.2 illustrates a spanning tree corresponding to a cycle free solution.

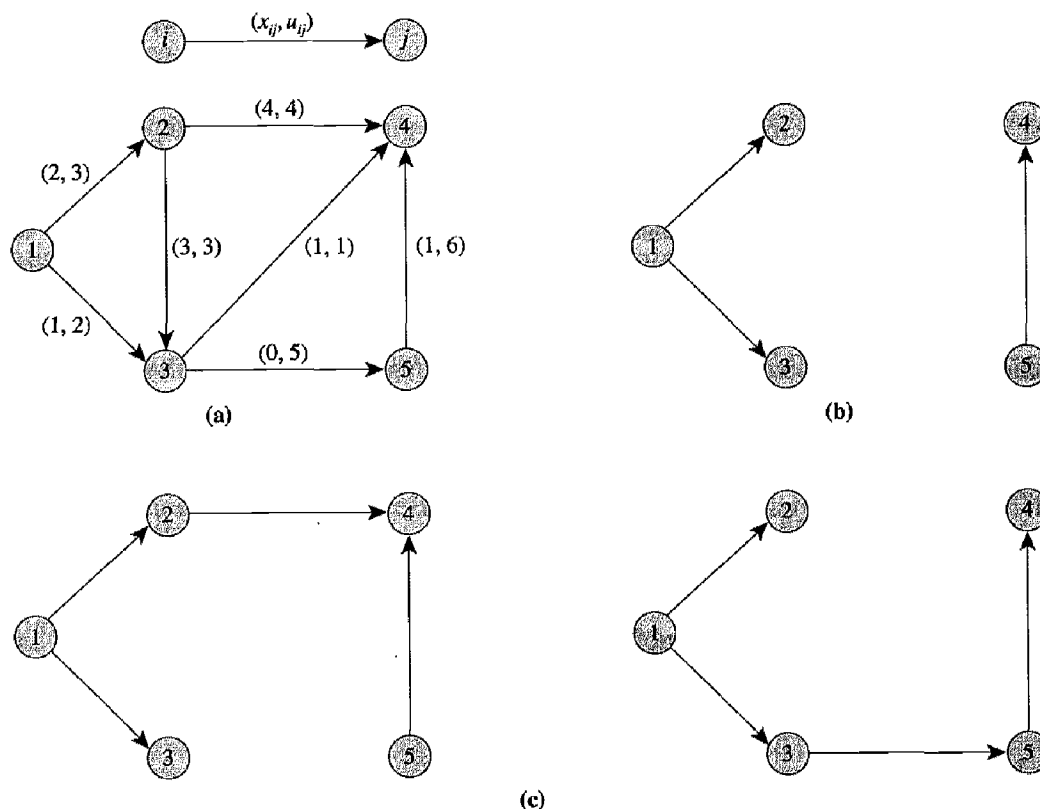


Figure 11.2 Converting a cycle free solution into a spanning tree solution: (a) example network; (b) set of free arcs; (c) 2 spanning tree solutions.

The solution in Figure 11.2(a) is cycle free. Figure 11.2(b) represents the set of free arcs, and Figure 11.2(c) shows two spanning tree solutions corresponding to the cycle free solution. As shown by this example, it might be possible (and often is) to complete the set of free arcs into a spanning tree in several ways. Adding the arc (3, 4) instead of the arc (2, 4) or (3, 5) would produce yet another spanning tree solution. Therefore, a given cycle free solution can correspond to several spanning trees. Nevertheless, since we assume that the underlying network is connected, we can always add some restricted arcs to the free arcs of a cycle free solution to produce a spanning tree, so we have established the following fundamental result:

Theorem 11.2 (Spanning Tree Property). *If the objective function of a minimum cost flow problem is bounded from below over the feasible region, the problem always has an optimal spanning tree solution.* ♦

A spanning tree solution partitions the arc set A into three subsets: (1) T , the arcs in the spanning tree; (2) L , the nontree arcs whose flow is restricted to value zero; and (3) U , the nontree arcs whose flow is restricted in value to the arcs' flow capacities. We refer to the triple (T, L, U) as a *spanning tree structure*.

Just as we can associate a spanning tree structure with a spanning tree solution, we can also obtain a unique spanning tree solution corresponding to a given spanning tree structure (T, L, U) . To do so, we set $x_{ij} = 0$ for all arcs $(i, j) \in L$, $x_{ij} = u_{ij}$ for all arcs $(i, j) \in U$, and then solve the mass balance equations to determine the flow values for arcs in T . In Section 11.4 we show that the flows on the spanning tree arcs are unique. We say that a spanning tree structure is *feasible* if its associated spanning tree solution satisfies all of the arcs' flow bounds. In the special case in which every tree arc in a spanning tree solution is a free arc, we say that the spanning tree is *nondegenerate*; otherwise, we refer to it as a *degenerate* spanning tree. We refer to a spanning tree structure as *optimal* if its associated spanning tree solution is an optimal solution of the minimum cost flow problem. The following theorem states a sufficient condition for a spanning tree structure to be an optimal structure. As shown by our discussion in previous chapters, the reduced costs defined as $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$ are useful in characterizing optimal solutions to minimum cost flow problems.

Theorem 11.3 (Minimum Cost Flow Optimality Conditions). *A spanning tree structure (T, L, U) is an optimal spanning tree structure of the minimum cost flow problem if it is feasible and for some choice of node potentials π , the arc reduced costs c_{ij}^π satisfy the following conditions:*

$$(a) \quad c_{ij}^\pi = 0 \text{ for all } (i, j) \in T. \quad (11.1a)$$

$$(b) \quad c_{ij}^\pi \geq 0 \text{ for all } (i, j) \in L. \quad (11.1b)$$

$$(c) \quad c_{ij}^\pi \leq 0 \text{ for all } (i, j) \in U. \quad (11.1c)$$

Proof. Let x^* be the solution associated with the spanning tree structure (T, L, U) . We know that some set of node potentials π , together with the spanning tree structure (T, L, U) , satisfies (11.1).

We need to show that x^* is an optimal solution of the minimum cost flow

problem. In Section 2.4 we showed that minimizing $\sum_{(i,j) \in A} c_{ij}x_{ij}$ is equivalent to minimizing $\sum_{(i,j) \in A} c_{ij}^{\pi}x_{ij}$. The conditions stated in (11.1) imply that for the given node potential π , minimizing $\sum_{(i,j) \in A} c_{ij}^{\pi}x_{ij}$ is equivalent to minimizing the following expression:

$$\text{Minimize } \sum_{(i,j) \in L} c_{ij}^{\pi}x_{ij} - \sum_{(i,j) \in U} |c_{ij}^{\pi}| x_{ij}. \quad (11.2)$$

The definition of the solution x^* implies that for any arbitrary solution x , $x_{ij} \geq x_{ij}^*$ for all $(i, j) \in L$ and $x_{ij} \leq x_{ij}^*$ for all $(i, j) \in U$. The expression (11.2) implies that the objective function value of the solution x will be greater than or equal to that of x^* . ♦

These optimality conditions have a nice economic interpretation. As we shall see later in Section 11.4, if $\pi(1) = 0$, the equations in (11.1a) imply that $-\pi(k)$ denotes the length of the tree path from node 1 to node k . The reduced cost $c_{ij}^{\pi} = c_{ij} - \pi(i) + \pi(j)$ for a nontree arc $(i, j) \in L$ denotes the change in the cost of the flow that we realize by sending 1 unit of flow through the tree path from node 1 to node i through the arc (i, j) , and then back to node 1 along the tree path from node j to node 1. The condition (11.1b) implies that this circulation of flow is not profitable (i.e., does not decrease cost) for any nontree arc in L . The condition (11.1c) has a similar interpretation.

The network simplex algorithm maintains a feasible spanning tree structure and moves from one spanning tree structure to another until it finds an optimal structure. At each iteration, the algorithm adds one arc to the spanning tree in place of one of its current arcs. The entering arc is a nontree arc violating its optimality condition. The algorithm (1) adds this arc to the spanning tree, creating a negative cycle (which might have zero residual capacity), (2) sends the maximum possible flow in this cycle until the flow on at least one arc in the cycle reaches its lower or upper bound, and (3) drops an arc whose flow has reached its lower or upper bound, giving us a new spanning tree structure. Because of its relationship to the primal simplex algorithm for the linear programming problem (see Appendix C), this operation of moving from one spanning tree structure to another is known as a *pivot operation*, and the two spanning trees structures obtained in consecutive iterations are called *adjacent spanning tree structures*. In Section 11.5 we give a detailed description of this algorithm.

11.3 MAINTAINING A SPANNING TREE STRUCTURE

Since the network simplex algorithm generates a sequence of spanning tree solutions, to implement the algorithm effectively, we need to be able to represent spanning trees conveniently in a computer so that the algorithm can perform its basic operations efficiently and can update the representation quickly when it changes the spanning tree. Over the years, researchers have suggested several procedures for maintaining and manipulating a spanning tree structure. In this section we describe one of the more popular representations.

We consider the tree as “hanging” from a specially designated node, called the *root*. Throughout this chapter we assume that node 1 is the root node. Figure

11.3 gives an example of a tree. We associate three indices with each node i in the tree: a predecessor index, $pred(i)$, a depth index $depth(i)$, and a thread index, $thread(i)$.

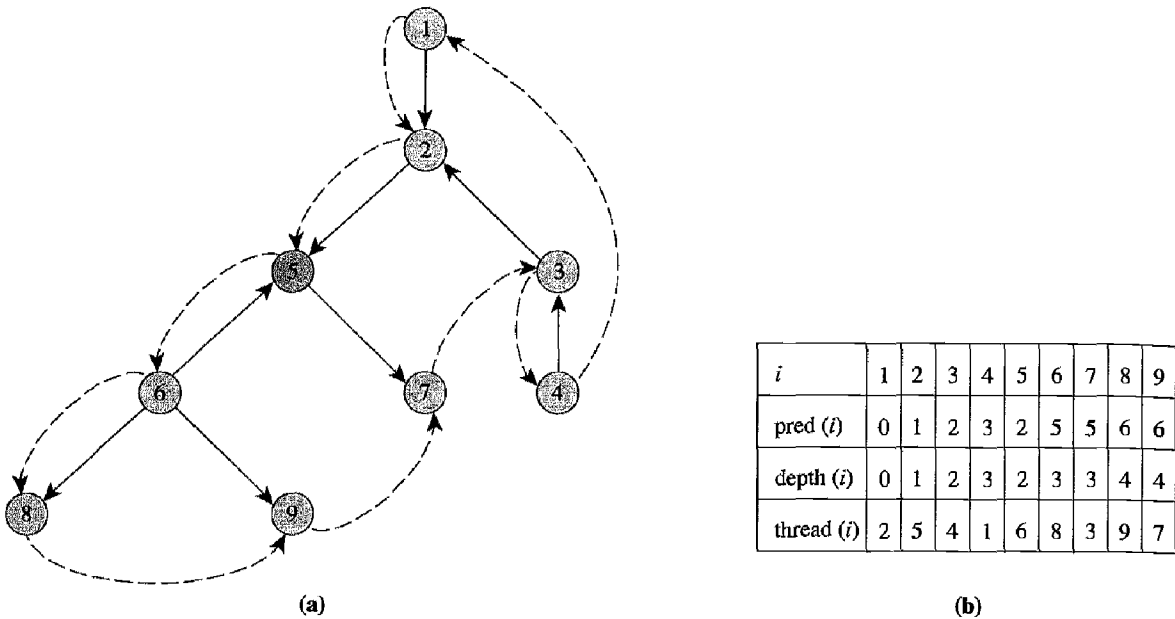


Figure 11.3 Example of a tree indices: (a) rooted tree; (b) corresponding tree indices.

Predecessor index. Each node i has a unique path connecting it to the root. The index $pred(i)$ stores the first node in that path (other than node i). For example, the path 9–6–5–2–1 connects node 9 to the root; therefore, $pred(9) = 6$. By convention, we set the predecessor node of the root node, node 1, equal to zero. Figure 11.3 specifies these indices for the other nodes. Observe that by iteratively using the predecessor indices, we can enumerate the path from any node to the root.

A node j is called a *successor* of node i if $pred(j) = i$. For example, node 5 has two successors: nodes 6 and 7. A *leaf node* is a node with no successors. In Figure 11.3, nodes 4, 7, 8, and 9 are leaf nodes. The *descendants* of a node i are the node i itself, its successors, successors of its successors, and so on. For example, in Figure 11.3, the elements of node set $\{5, 6, 7, 8, 9\}$ are the descendants of node 5.

Depth index. We observed earlier that each node i has a unique path connecting it to the root. The index $depth(i)$ stores the number of arcs in that path. For example, since the path 9–6–5–2–1 connects node 9 to the root, $depth(9) = 4$. Figure 11.3 gives depth indices for all of the nodes in the network.

Thread index. The thread indices define a traversal of a tree, that is, a sequence of nodes that walks or threads its way through the nodes of a tree, starting at the root node, and visiting nodes in a “top-to-bottom” order, and finally returning to the root. We can find thread indices by performing a depth-first search of the tree

as described in Section 3.4 and setting the thread of a node to be the node in the depth-first search encountered just after the node itself. For our example, the depth-first traversal would read 1–2–5–6–8–9–7–3–4–1, so $\text{thread}(1) = 2$, $\text{thread}(2) = 5$, $\text{thread}(5) = 6$, and so on (see the dashed lines in Figure 11.3).

The thread indices provide a particularly convenient means for visiting (or finding) all descendants of a node i . We simply follow the thread starting at that node and record the nodes visited, until the depth of the visited node becomes at least as large as that of node i . For example, starting at node 5, we visit nodes 6, 8, 9, and 7 in order, which are the descendants of node 5 and then visit node 3. Since the depth of node 3 equals that of node 5, we know that we have left the “descendant tree” lying below node 5. We shall see later that finding the descendant tree of a node efficiently is an important step in developing an efficient implementation of the network simplex algorithm.

In the next section we show how the tree indices permit us to compute the feasible solution and the set of node potentials associated with a tree.

11.4 COMPUTING NODE POTENTIALS AND FLOWS

As we noted in Section 11.2, as the network simplex algorithm moves from one spanning tree to the next, it always maintains the condition that the reduced cost of every arc (i, j) in the current spanning tree is zero (i.e. $c_{ij}^{\pi} = 0$). Given the current spanning tree structure (T, L, U) , the method first determines values for the node potentials π that will satisfy this condition for the tree arcs. In this section we show how to find these values of the node potentials.

Note that we can set the value of one node potential arbitrarily because adding a constant k to each node potential does not alter the reduced cost of any arc; that is, for any constant k , $c_{ij}^{\pi} = c_{ij} - \pi(i) + \pi(j) = c_{ij} - [\pi(i) + k] + [\pi(j) + k]$. So for convenience, we henceforth assume that $\pi(1) = 0$. We compute the remaining node potentials using the fact that the reduced cost of every spanning tree arc is zero; that is,

$$c_{ij}^{\pi} = c_{ij} - \pi(i) + \pi(j) = 0 \quad \text{for every arc } (i, j) \in T. \quad (11.3)$$

In equation (11.3), if we know one of the node potentials $\pi(i)$ or $\pi(j)$, we can easily compute the other one. Consequently, the basic idea in the procedure is to start at node 1 and fan out along the tree arcs using the thread indices to compute other node potentials. By traversing the nodes using the thread indices, we ensure that whenever the procedure visits a node k , it has already evaluated the potential of its predecessor, so it can compute $\pi(k)$ using (11.3). Figure 11.4 gives a formal statement of the procedure *compute-potentials*.

The numerical example shown in Figure 11.5 illustrates the procedure. We first set $\pi(1) = 0$. The thread of node 1 is 2, so we next examine node 2. Since arc $(1, 2)$ connects node 2 to its predecessor, using (11.3) we find that $\pi(2) = \pi(1) - c_{12} = -5$. We next examine node 5, which is connected to its parent by arc $(5, 2)$. Using (11.3) we obtain $\pi(5) = \pi(2) + c_{52} = -5 + 2 = -3$. In the same fashion we compute the rest of the node potentials; the numbers shown next to each node in Figure 11.5 specify these values.

```

procedure compute-potentials;
begin
   $\pi(1) := 0$ ;
   $j := \text{thread}(1)$ ;
  while  $j \neq 1$  do
    begin
       $i := \text{pred}(j)$ ;
      if  $(i, j) \in A$  then  $\pi(j) := \pi(i) - c_{ij}$ ;
      if  $(j, i) \in A$  then  $\pi(j) := \pi(i) + c_{ji}$ ;
       $j := \text{thread}(j)$ ;
    end;
  end;

```

Figure 11.4 Procedure *compute-potentials*.

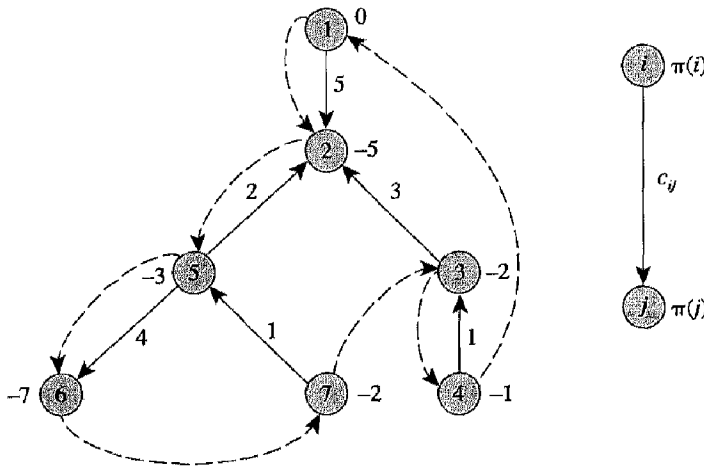


Figure 11.5 Computing node potentials for a spanning tree.

Let P be the tree path in T from the root node 1 to some node k . Moreover, let \bar{P} and \underline{P} , respectively, denote the sets of forward and backward arcs in P . Now let us examine arcs in P starting at node 1. The procedure *compute-potentials* implies that $\pi(j) = \pi(i) - c_{ij}$ whenever arc (i, j) is a forward arc in the path, and that $\pi(j) = \pi(i) + c_{ji}$ whenever arc (j, i) is a backward arc in the path. This observation implies that $\pi(k) = \pi(k) - \pi(1) = -\sum_{(i,j) \in \bar{P}} c_{ij} + \sum_{(i,j) \in \underline{P}} c_{ji}$. In other words, $\pi(k)$ is the negative of the cost of sending 1 unit of flow from node 1 to node k along the tree path. Alternatively, $\pi(k)$ is the cost of sending 1 unit of flow from node k to node 1 along the tree path. The procedure *compute-potentials* requires $O(1)$ time per iteration and performs $(n - 1)$ iterations to evaluate the node potential of each node. Therefore, the procedure runs in $O(n)$ time.

One important consequence of the procedure *compute-potentials* is that the minimum cost flow problem always has integer optimal node potentials whenever all the arc costs are integer. To see this result, recall from Theorem 11.2 that the minimum cost flow problem always has an optimal spanning tree solution. The potentials associated with this tree constitute optimal node potentials, which we can determine using the procedure *compute-potentials*. The description of the procedure *compute-potentials* implies that if all arc costs are integer, node potentials are integer as well (because the procedure performs only additions and subtractions). We refer to this integrality property of optimal node potentials as the *dual integrality property*.

since node potentials are the dual linear programming variables associated with the minimum cost flow problem.

Theorem 11.4 (Dual Integrality Property). *If all arc costs are integer, the minimum cost flow problem always has optimal integer node potentials.* ♦

Computing Arc Flows

We next consider the problem of determining the flows on the tree arcs of a given spanning tree structure. To ease our discussion, for the moment let us first consider the *uncapacitated* version of the minimum cost flow problem. We can then assume that all nontree arcs carry zero flow.

If we delete a tree arc, say arc (i, j) , from the spanning tree, the tree decomposes into two subtrees. Let T_1 be the subtree containing node i and let T_2 be the subtree containing node j . Note that $\sum_{k \in T_1} b(k)$ denotes the cumulative supply/demand of nodes in T_1 [which must be equal to $-\sum_{k \in T_2} b(k)$ because $\sum_{k \in T_1} b(k) + \sum_{k \in T_2} b(k) = 0$]. In the spanning tree, arc (i, j) is the only arc that connects the subtree T_1 to the subtree T_2 , so it must carry $\sum_{k \in T_1} b(k)$ units of flow, for this is the only way to satisfy the mass balance constraints. For example, in Figure 11.6, if we delete arc $(1, 2)$ from the tree, then $T_1 = \{1, 3, 6, 7\}$, $T_2 = \{2, 4, 5\}$, and $\sum_{k \in T_1} b(k) = 10$. Consequently, arc $(1, 2)$ carries 10 units of flow.

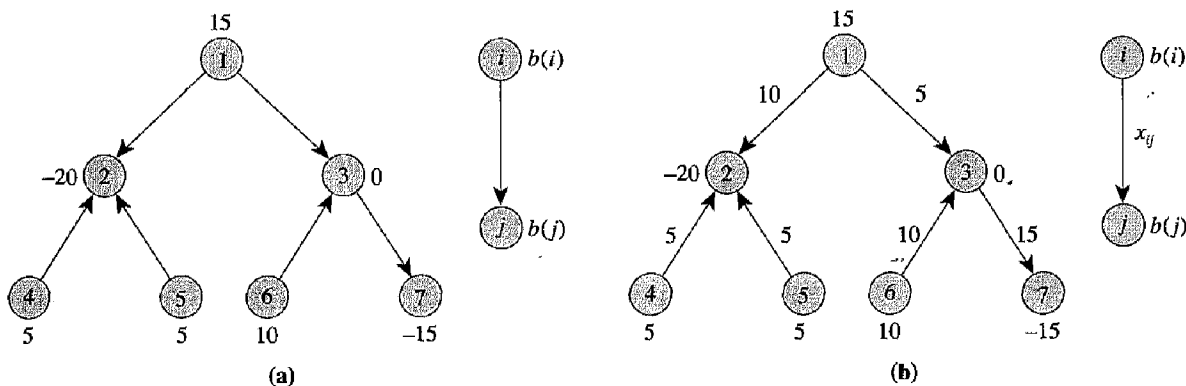


Figure 11.6 Computing flows for a spanning tree.

Using this observation we can devise an efficient method for computing the flows on all the tree arcs. Suppose that (i, j) is a tree arc and that node j is a leaf node [the treatment of the case when (i, j) is a tree arc and node i is a leaf node is similar]. Our observations imply that arc (i, j) must carry $-b(j)$ units of flow. For our example, arc $(3, 7)$ must carry 15 units of flow to satisfy the demand of node 7. Setting the flow on this arc to this value has an effect on the mass balance of its incident nodes: we must subtract 15 units from $b(3)$ and add 15 units to $b(7)$ [which reduces $b(7)$ to zero]. Having determined x_{37} , we can delete arc $(3, 7)$ from the tree and repeat the method on the smaller tree. Notice that we can identify a leaf node in every iteration because every tree has at least two leaf nodes (see Exercise 2.13). Figure 11.7 gives a formal description of this procedure.

```

procedure compute-flows;
begin
   $b'(i) := b(i)$ , for all  $i \in N$ ;
  for each  $(i, j) \in L$  do set  $x_{ij} := 0$ ;
   $T' := T$ ;
  while  $T' \neq \{1\}$  do
    begin
      select a leaf node  $j$  (other than node 1) in the subtree  $T'$ ;
       $i := \text{pred}(j)$ ;
      if  $(i, j) \in T'$  then  $x_{ij} := -b'(j)$ 
      else  $x_{ij} := b'(j)$ ;
      add  $b'(j)$  to  $b'(i)$ ;
      delete node  $j$  and the arc incident to it from  $T'$ ;
    end;
  end;

```

Figure 11.7 Procedure *compute-flows*.

This method for computing the flow values assumes that the minimum cost flow problem is uncapacitated. For the capacitated version of the problem, we add the following statement immediately after the first statement [i.e., $b'(i) := b(i)$ for all $i \in N$] in the procedure *compute-flows*. We leave the justification of this modification as an exercise (see Exercise 11.19).

```

for each  $(i, j) \in U$  do
  set  $x_{ij} := u_{ij}$ , subtract  $u_{ij}$  from  $b'(i)$  and add  $u_{ij}$  to  $b'(j)$ ;

```

The running time of the procedure *compute-flows* is easy to determine. Clearly, the initialization of flows and modification of supplies/demands $b(i)$ and $b(j)$ for arcs (i, j) in U requires $O(m)$ time. If we set aside the time to select leaf nodes of T , then each iteration requires $O(1)$ time, resulting in a total of $O(n)$ time. One way of identifying leaf nodes in T is to select nodes in the reverse order of the thread indices. Note that in the thread traversal, each node appears prior to its descendants (see Property 3.4). We identify the reverse thread traversal of the nodes by examining the nodes in the order dictated by the thread indices, putting all the nodes into a stack in the order of their appearance and then taking them out from the top of the stack one at a time. Therefore, the reverse thread traversal examines each node only after it has examined all of the node's descendants. We have thus established that for the uncapacitated minimum cost flow problem, the procedure *compute-flows* runs in $O(m)$ time. For the capacitated version of the problem, the procedure also requires $O(m)$ time.

We can use the procedure *compute-flows* to obtain an alternative proof of the (primal) integrality property that we stated in Theorem 9.10. Recall from Theorem 11.2 that the minimum cost flow problem always has an optimal spanning tree solution. The flow associated with this tree is an optimal flow and we can determine it using the procedure *compute-flows*. The description of the procedure *compute-flows* implies that if the capacities of all the arcs and the supplies/demands of all the nodes are integer, arc flows are integer as well (because the procedure performs only additions and subtractions). We state this result again because of its importance in network flow theory.

Theorem 11.5 (Primal Integrality Property). *If capacities of all the arcs and supplies/demands of all the nodes are integer, the minimum cost flow problem always has an integer optimal flow.* ♦

In closing this section we observe that every spanning tree structure (T, L, U) defines a unique flow x . If this flow satisfies the flow bounds $0 \leq x_{ij} \leq u_{ij}$ for every arc $(i, j) \in A$, the spanning tree structure is feasible; otherwise, it is infeasible. We refer to the spanning tree T as *degenerate* if $x_{ij} = 0$ or $x_{ij} = u_{ij}$ for some arc $(i, j) \in T$, and *nondegenerate* otherwise. In a nondegenerate spanning tree, $0 < x_{ij} < u_{ij}$ for every tree arc (i, j) .

11.5 NETWORK SIMPLEX ALGORITHM

The network simplex algorithm maintains a feasible spanning tree structure at each iteration and successively transforms it into an improved spanning tree structure until it becomes optimal. The algorithmic description in Figure 11.8 specifies the essential steps of the method.

```

algorithm network simplex;
begin
    determine an initial feasible tree structure  $(T, L, U)$ ;
    let  $x$  be the flow and  $\pi$  be the node potentials associated with this tree structure;
    while some nontree arc violates the optimality conditions do
        begin
            select an entering arc  $(k, l)$  violating its optimality condition;
            add arc  $(k, l)$  to the tree and determine the leaving arc  $(p, q)$ ;
            perform a tree update and update the solutions  $x$  and  $\pi$ ;
        end;
    end;

```

Figure 11.8 Network simplex algorithm.

In the following discussion we describe in greater detail how the network simplex algorithm uses tree indices to perform these various steps. This discussion highlights the value of the tree indices in designing an efficient implementation of the algorithm.

Obtaining an Initial Spanning Tree Structure

Our connectedness assumption (i.e., Assumption 9.4 in Section 9.1) provides one way of obtaining an initial spanning tree structure. We have assumed that for every node $j \in N - \{1\}$, the network contains arcs $(1, j)$ and $(j, 1)$, with sufficiently large costs and capacities. We construct the initial tree T as follows. We examine each node j , other than node 1, one by one. If $b(j) \geq 0$, we include arc $(1, j)$ in T with a flow value of $b(j)$. If $b(j) < 0$, we include arc $(j, 1)$ in T with a flow value of $-b(j)$. The set L consists of the remaining arcs, and the set U is empty. As shown in Section 11.4, we can easily compute the node potentials for this tree using the equations $c_{ij} - \pi(i) + \pi(j) = 0$ for all $(i, j) \in T$. Recall that we set $\pi(1) = 0$.

If the network does not contain the arcs $(1, j)$ and $(j, 1)$ for each node $j \in$

$N - \{1\}$ (or, we do not wish to add these arcs for some reason), we could construct an initial spanning tree structure by first establishing a feasible flow in the network by solving a maximum flow problem (as described in Application 6.1), and then by converting this solution into a spanning tree solution using the method described in Section 11.2.

Optimality Testing and the Entering Arc

Let (T, L, U) be a feasible spanning tree structure of the minimum cost flow problem, and let π be the corresponding node potentials. To determine whether the spanning tree structure is optimal, we check to see whether the spanning tree structure satisfies the following conditions:

$$c_{ij}^{\pi} \geq 0 \text{ for every arc } (i, j) \in L,$$

$$c_{ij}^{\pi} \leq 0 \text{ for every arc } (i, j) \in U.$$

If the spanning tree structure satisfies these conditions, it is optimal and the algorithm terminates. Otherwise, the algorithm selects a nontree arc violating the optimality condition to be introduced into the tree. Two types of arcs are *eligible* to enter the tree:

1. Any arc $(i, j) \in L$ with $c_{ij}^{\pi} < 0$
2. Any arc $(i, j) \in U$ with $c_{ij}^{\pi} > 0$

For any eligible arc (i, j) , we refer to $|c_{ij}^{\pi}|$ as its *violation*. The network simplex algorithm can select *any* eligible arc to enter the tree and still would terminate finitely (with some provisions for dealing with degeneracy, as discussed in Section 11.6). However, different rules for selecting the entering arc produce algorithms with different empirical and theoretical behavior. Many different rules, called *pivot rules*, are possible for choosing the entering arc. The following rules are most widely adopted.

Dantzig's pivot rule. This rule was suggested by George B. Dantzig, the father of linear programming. At each iteration this rule selects an arc with the maximum violation to enter the tree. The motivation for this rule is that the arc with the maximum violation causes the maximum decrease in the objective function per unit change in the value of flow on the selected arc, and hence the introduction of this arc into the spanning tree would cause the maximum decrease per pivot if the average increase in the value of the selected arc were the same for all arcs. Computational results confirm that this choice of the entering arc tends to produce relatively large decreases in the objective function per iteration and, as a result, the algorithm performs fewer iterations than other choices for the pivot rule. However, this rule does have a major drawback: The algorithm must consider every nontree arc to identify the arc with the maximum violation and doing so is very time consuming. Therefore, even though this algorithm generally performs fewer iterations than other implementations, the running time of the algorithm is not attractive.

First eligible arc pivot rule. To implement this rule, we scan the arc list sequentially and select the first eligible arc to enter the tree. In a popular version of this rule, we examine the arc list in a wraparound fashion. For example, in an iteration if we find that the fifth arc in the arc list is the first eligible arc, then in the next iteration we start scanning the arc list from the sixth arc. If we reach the end of the arc list while we are performing some iteration, we continue by examining the arc list from the beginning. One nice feature of this pivot rule is that it quickly identifies the entering arc. The pivot rule does have a counterbalancing drawback: with it, the algorithm generally performs more iterations than it would with other pivot rules because each pivot operation produces a relatively small decrease in the objective function value. The overall effect of this pivot rule on the running time of the algorithm is not very attractive, although the rule does produce a more efficient implementation than Dantzig's pivot rule.

Dantzig's pivot rule and the first pivot rule represent two extreme choices of a pivot rule. The *candidate list pivot rule*, which we discuss next, strikes an effective compromise between these two extremes and has proven to be one of the most successful pivot rules in practice. This rule also offers sufficient flexibility for fine tuning to special circumstances.

Candidate list pivot rule. When implemented with this rule, the algorithm selects the entering arc using a two-phase procedure consisting of *major iterations* and *minor iterations*. In a major iteration we construct a *candidate list* of eligible arcs. Having constructed this list, we then perform a number of minor iterations; in each of these iterations, we select an eligible arc from the candidate list with the maximum violation.

In a major iteration we construct the candidate list as follows. We first examine arcs emanating from node 1 and add eligible arcs to the candidate list. We repeat this process for nodes 2, 3, . . . , until either the list has reached its maximum allowable size or we have examined all the nodes. The next major iteration begins with the node where the previous major iteration ended and examines nodes in a wraparound fashion.

Once the algorithm has formed the candidate list in a major iteration, it performs a number of minor iterations. In a minor iteration, the algorithm scans all the arcs in the candidate list and selects an arc with the maximum violation to enter the tree. As we scan the arcs, we update the candidate list by removing those arcs that are no longer eligible (due to changes in the node potentials). Once the candidate list becomes empty or we have reached a specified limit on the number of minor iterations to be performed within each major iteration, we rebuild the candidate list by performing another major iteration.

Notice that the candidate list approach offers considerable flexibility for fine tuning to special problem classes. By setting the maximum allowable size of the candidate list appropriately and by specifying the number of minor iterations to be performed within a major iteration, we can obtain numerous different pivot rules. In fact, Dantzig's pivot rule and the first eligible pivot rule are special cases of the candidate list pivot rule (see Exercise 11.20).

In the preceding discussion, we described several important pivot rules. In the reference notes, we supply references for other pivot rules. Our next topic of study

is deciding how to choose the arc that leaves the spanning tree structure at each step of the network simplex algorithm.

Leaving Arc

Suppose that we select arc (k, l) as the entering arc. The addition of this arc to the tree T creates exactly one cycle W , which we refer to as the *pivot cycle*. The pivot cycle consists of the unique path in the tree T from node k to node l , together with arc (k, l) . We define the orientation of the cycle W as the same as that of (k, l) if $(k, l) \in L$ and opposite the orientation of (k, l) if $(k, l) \in U$. Let \overline{W} and \underline{W} denote the sets of *forward arcs* (i.e., those along the orientation of W) and *backward arcs* (those opposite to the orientation of W) in the pivot cycle. Sending additional flow around the pivot cycle W in the direction of its orientation strictly decreases the cost of the current solution at the per unit rate of $|c_{kl}|$. We augment the flow as much as possible until one of the arcs in the pivot cycle reaches its lower or upper bound. Notice that augmenting flow along W increases the flow on forward arcs and decreases the flow on backward arcs. Consequently, the maximum flow change δ_{ij} on an arc $(i, j) \in W$ that satisfies the flow bound constraints is

$$\delta_{ij} = \begin{cases} u_{ij} - x_{ij} & \text{if } (i, j) \in \overline{W} \\ x_{ij} & \text{if } (i, j) \in \underline{W} \end{cases}$$

To maintain feasibility, we can augment $\delta = \min\{\delta_{ij} : (i, j) \in W\}$ units of flow along W . We refer to any arc $(i, j) \in W$ that defines δ (i.e., for which $\delta = \delta_{ij}$) as a *blocking arc*. We then augment δ units of flow and select an arc (p, q) with $\delta_{pq} = \delta$ as the leaving arc, breaking ties arbitrarily. We say that a pivot iteration is a *nondegenerate iteration* if $\delta > 0$ and is a *degenerate iteration* if $\delta = 0$. A degenerate iteration occurs only if T is a degenerate spanning tree. Observe that if two arcs tie while determining the value of δ , the next spanning tree will be degenerate.

The crucial step in identifying the leaving arc is to identify the pivot cycle. If $P(i)$ denotes the unique path in the tree from any node i to the root node, this cycle consists of the arcs $\{(k, l)\} \cup P(k) \cup P(l) - (P(k) \cap P(l))$. In other words, W consists of the arc (k, l) and the disjoint portions of $P(k)$ and $P(l)$. Using the predecessor indices alone permits us to identify the cycle W as follows. First, we designate all the nodes in the network as unmarked. We then start at node k and, using the predecessor indices, trace the path from this node to the root and mark all the nodes in this path. Next we start at node l and trace the predecessor indices until we encounter a marked node, say w . The node w is the first common ancestor of nodes k and l ; we refer to it as the *apex* of cycle W . The cycle W contains the portions of the paths $P(k)$ and $P(l)$ up to node w , together with the arc (k, l) . This method identifies the cycle W in $O(n)$ time and so is efficient. However, it has the drawback of backtracking along those arcs of $P(k)$ that are not in W . If the pivot cycle lies “deep in the tree,” far from its root, then tracing the nodes back to the root will be inefficient. Ideally, we would like to identify the cycle W in time proportional to $|W|$. The simultaneous use of depth and predecessor indices, as indicated in Figure 11.9, permits us to achieve this goal.

This method scans the arcs in the pivot cycle W twice. During the first scan, we identify the apex of the cycle and also identify the maximum possible flow that

```

procedure identify-cycle;
begin
   $i := k$  and  $j := l$ ;
  while  $i \neq j$  do
    begin
      if  $\text{depth}(i) > \text{depth}(j)$  then  $i := \text{pred}(i)$ 
      else if  $\text{depth}(j) > \text{depth}(i)$  then  $j := \text{pred}(j)$ 
      else  $i := \text{pred}(i)$  and  $j := \text{pred}(j)$ ;
    end;
  end;

```

Figure 11.9 Procedure for identifying the pivot cycle.

can be augmented along W . In the second scan, we augment the flow. The entire flow change operation requires $O(n)$ time in the worst case, but typically it examines only a small subset of nodes (and arcs).

Updating the Tree

When the network simplex algorithm has determined a leaving arc (p, q) for a given entering arc (k, l) , it updates the tree structure. If the leaving arc is the same as the entering arc, which would happen when $\delta = \delta_{kl} = u_{kl}$, the tree does not change. In this instance the arc (k, l) merely moves from the set \mathbf{L} to the set \mathbf{U} , or vice versa. If the leaving arc differs from the entering arc, the algorithm must perform more extensive changes. In this instance the arc (p, q) becomes a nontree arc at its lower or upper bound, depending on whether (in the updated flow) $x_{pq} = 0$ or $x_{pq} = u_{pq}$. Adding arc (k, l) to the current spanning tree and deleting arc (p, q) creates a new spanning tree.

For the new spanning tree, the node potentials also change; we can update them as follows. The deletion of the arc (p, q) from the previous tree partitions the set of nodes into two subtrees, one, T_1 , containing the root node, and the other, T_2 , not containing the root node. Note that the subtree T_2 hangs from node p or node q . The arc (k, l) has one endpoint in T_1 and the other in T_2 . As is easy to verify, the conditions $\pi(1) = 0$ and $c_{ij} - \pi(i) + \pi(j) = 0$ for all arcs in the new tree imply that the potentials of nodes in the subtree T_1 remain unchanged, and the potentials of nodes in the subtree T_2 change by a constant amount. If $k \in T_1$ and $l \in T_2$, all the node potentials in T_2 increase by $-c_{kl}^{\pi}$; if $l \in T_1$ and $k \in T_2$, they increase by the amount c_{kl}^{π} . Using the thread and depth indices, the method described in Figure 11.10 updates the node potentials quickly.

```

procedure update-potentials;
begin
  if  $q \in T_2$  then  $y := q$  else  $y := p$ ;
  if  $k \in T_1$  then  $\text{change} := -c_{kl}^{\pi}$  else  $\text{change} := c_{kl}^{\pi}$ ;
   $\pi(y) := \pi(y) + \text{change}$ ;
   $z := \text{thread}(y)$ ;
  while  $\text{depth}(z) > \text{depth}(y)$  do
    begin
       $\pi(z) := \pi(z) + \text{change}$ ;
       $z := \text{thread}(z)$ ;
    end;
  end;

```

Figure 11.10 Updating node potentials in a pivot operation.

The final step in the updating of the tree is to recompute the various tree indices. This step is rather involved and we refer the reader to the references given in reference notes for the details. We do point out, however, that it is possible to update the tree indices in $O(n)$ time. In fact, the time required to update the tree indices is $O(|W| + \min\{|T_1|, |T_2|\})$, which is typically much less than n .

Termination

The network simplex algorithm, as just described, moves from one feasible spanning tree structure to another until it obtains a spanning tree structure that satisfies the optimality condition (11.1). If each pivot operation in the algorithm is nondegenerate, it is easy to show that the algorithm terminates finitely. Recall that $|c_{kl}^*|$ is the net decrease in the cost per unit flow sent around the pivot cycle W . After a nondegenerate pivot (for which $\delta > 0$), the cost of the new spanning tree structure is $\delta|c_{kl}^*|$ units less than the cost of the previous spanning tree structure. Since any network has a finite number of spanning tree structures and every spanning tree structure has a unique associated cost, the network simplex algorithm will encounter any spanning tree structure at most once and hence will terminate finitely. Degenerate pivots, however, pose a theoretical difficulty: The algorithm might not terminate finitely unless we perform pivots carefully. In the next section we discuss a special implementation, called the *strongly feasible spanning tree implementation*, that guarantees finite convergence of the network simplex algorithm even for problems that are degenerate.

We use the example in Figure 11.11(a) to illustrate the network simplex algorithm. Figure 11.11(b) shows a feasible spanning tree solution for the problem. For this solution, $T = \{(1, 2), (1, 3), (2, 4), (2, 5), (5, 6)\}$, $L = \{(2, 3), (5, 4)\}$, and $U = \{(3, 5), (4, 6)\}$. In this solution, arc $(3, 5)$ has a positive violation, which is 1 unit. We introduce this arc into the tree creating a cycle whose apex is node 1. Since arc $(3, 5)$ is at its upper bound, the orientation of the cycle is opposite to that of arc $(3, 5)$. The arcs $(1, 2)$ and $(2, 5)$ are forward arcs in the cycle and arcs $(3, 5)$ and $(1, 3)$ are backward arcs. The maximum increase in flow permitted by the arcs $(3, 5)$, $(1, 3)$, $(1, 2)$, and $(2, 5)$ is, respectively, 3, 3, 2, and 1 units. Consequently, $\delta = 1$ and we augment 1 unit of flow along the cycle. The augmentation increases the flow on arcs $(1, 2)$ and $(2, 5)$ by one unit and decreases the flow on arcs $(1, 3)$ and $(3, 5)$ by one unit. Arc $(2, 5)$ is the unique blocking arc and so we select it to leave the tree. Dropping arc $(2, 5)$ from the tree produces two subtrees: T_1 consisting of nodes 1, 2, 3, 4 and T_2 consisting of nodes 5 and 6. Introducing arc $(3, 5)$, we again obtain a spanning tree, as shown in Figure 11.11(c). Notice that in this spanning tree, the node potentials of nodes 5 and 6 are 1 unit less than that in the previous spanning tree.

In the feasible spanning tree solution shown in Figure 11.11(c), $L = \{(2, 3), (5, 4)\}$ and $U = \{(2, 5), (4, 6)\}$. In this solution, arc $(4, 6)$ is the only eligible arc: its violation equals 1 unit. Therefore, we introduce arc $(4, 6)$ into the tree. Figure 11.11(c) shows the resulting cycle and its orientation. We can augment 1 unit of additional flow along the orientation of this cycle. Sending this flow, we find that arc $(3, 5)$ is a blocking arc, so we drop this arc from the current spanning tree. Figure 11.11(d)

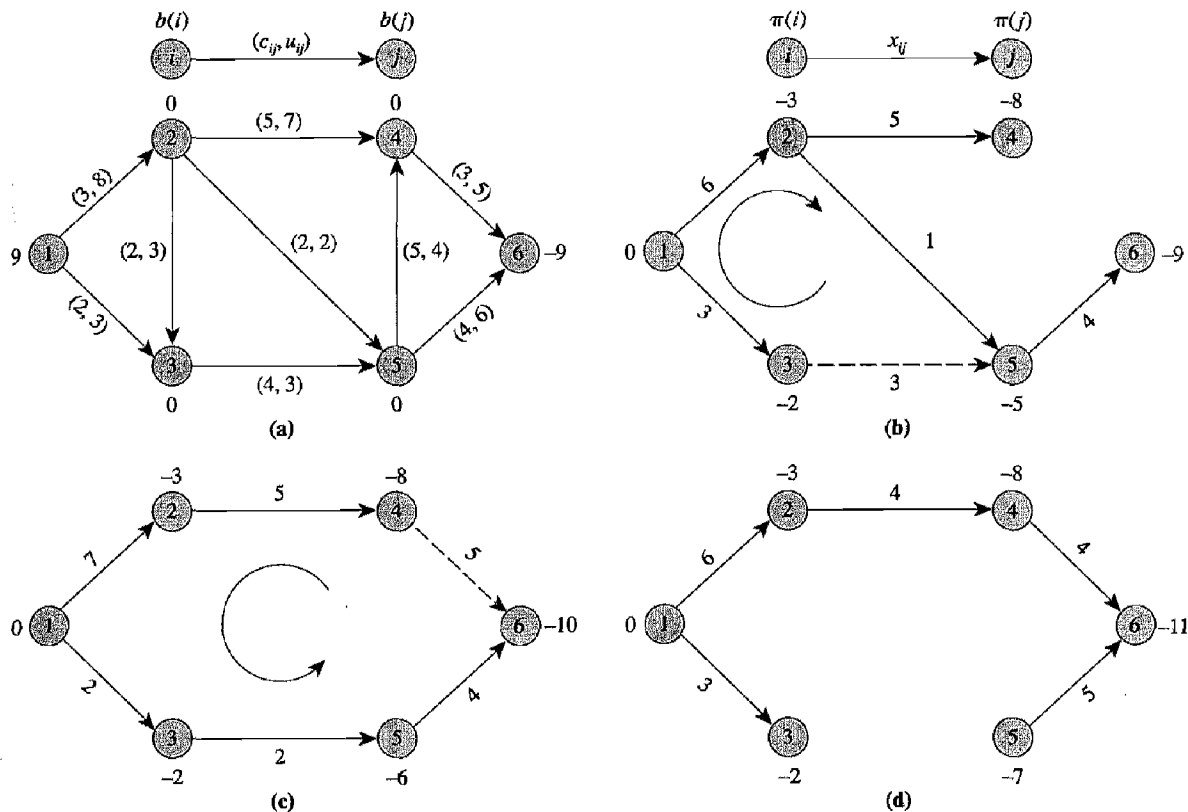


Figure 11.11 Numerical example for the network simplex algorithm.

shows the new spanning tree. As the reader can verify, this solution has no eligible arc, and thus the network simplex algorithm terminates with this solution.

11.6 STRONGLY FEASIBLE SPANNING TREES

The network simplex algorithm does not necessarily terminate in a finite number of iterations unless we impose some additional restriction on the choice of the entering and leaving arcs. Very small network examples show that a poor choice leads to *cycling* (i.e., an infinite repetitive sequence of degenerate pivots). Degeneracy in network problems is not only a theoretical issue, but also a practical one. Computational studies have shown that as many as 90% of the pivot operations in commonplace networks can be degenerate. As we show next, by maintaining a special type of spanning tree, called a *strongly feasible spanning tree*, the network simplex algorithm terminates finitely; moreover, it runs faster in practice as well.

Let (T, L, U) be a spanning tree structure for a minimum cost flow problem with integral data. As before, we conceive of a spanning tree as a tree hanging from the root node. The tree arcs are either *upward pointing* (toward the root) or are *downward pointing* (away from the root). We now state two alternate definitions of a strongly feasible spanning tree.

1. *Strongly feasible spanning tree.* A spanning tree T is *strongly feasible* if every tree arc with zero flow is upward pointing and every tree arc whose flow equals its capacity is downward pointing.
2. *Strongly feasible spanning tree.* A spanning tree T is *strongly feasible* if we can send a positive amount of flow from any node to the root along the tree path without violating any flow bound.

If a spanning tree T is strongly feasible, we also say that the spanning tree structure (T, L, U) is strongly feasible.

It is easy to show that the two definitions of the strongly feasible spanning trees are equivalent (see Exercise 11.24). Figure 11.12(a) gives an example of a strongly feasible spanning tree, and Figure 11.12(b) illustrates a feasible spanning tree that is not strongly feasible. The spanning tree shown in Figure 11.12(b) fails to be strongly feasible because arc $(3, 5)$ carries zero flow and is downward pointing. Observe that in this spanning tree, we cannot send any additional flow from nodes 5 and 7 to the root along the tree path.

To implement the network simplex algorithm so that it always maintains a strongly feasible spanning tree, we must first find an initial strongly feasible spanning tree. The method described in Section 11.5 for constructing the initial spanning tree structure always gives such a spanning tree. Note that a nondegenerate spanning tree is always strongly feasible; a degenerate spanning tree might or might not be strongly feasible. The network simplex algorithm creates a degenerate spanning tree from a nondegenerate spanning tree whenever two or more arcs are qualified as

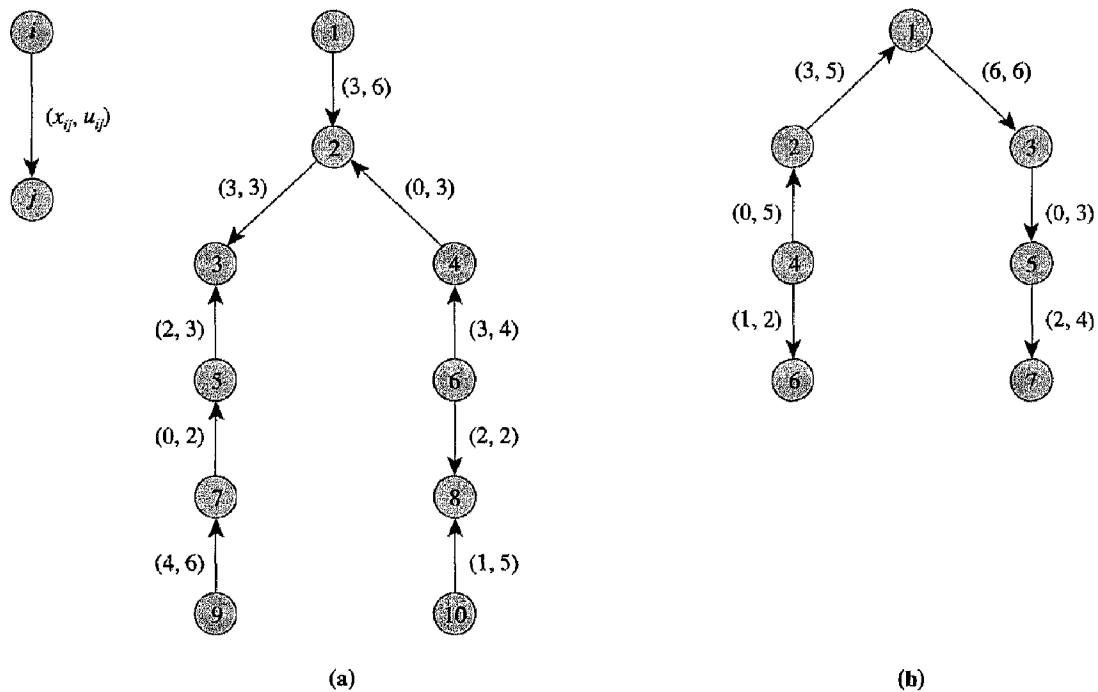


Figure 11.12 Feasible spanning trees: (a) strongly feasible; (b) nonstrongly feasible.

leaving arcs and we drop only one of these. Therefore, the algorithm needs to select the leaving arc carefully so that the next spanning tree is strongly feasible.

Suppose that we have a strongly feasible spanning tree and, during a pivot operation, arc (k, l) enters the spanning tree. We first consider the case when (k, l) is a nontree arc at its lower bound. Suppose that W is the pivot cycle formed by adding arc (k, l) to the spanning tree and that node w is the apex of the cycle W ; that is, w is the first common ancestor of nodes k and l . We define the orientation of the cycle W as compatible with that of arc (k, l) . After augmenting δ units of flow along the pivot cycle, the algorithm identifies the *blocking arcs* [i.e., those arcs (i, j) in the cycle that satisfy $\delta_{ij} = \delta$]. If the blocking arc is unique, we select it to leave the spanning tree. If the cycle contains more than one blocking arc, the next spanning tree will be degenerate (i.e., some tree arcs will be at their lower or upper bounds). In this case the algorithm selects the leaving arc in accordance with the following rule.

Leaving Arc Rule. *Select the leaving arc as the last blocking arc encountered in traversing the pivot cycle W along its orientation starting at the apex w .*

To illustrate the leaving arc rule, we consider a numerical example. Figure 11.13 shows a strongly feasible spanning tree for this example. Let $(9, 10)$ be the entering arc. The pivot cycle is $10-8-6-4-2-3-5-7-9-10$ and the apex is node 2. This pivot is degenerate because arcs $(2, 3)$ and $(7, 5)$ block any additional flow in the pivot cycle. Traversing the pivot cycle starting at node 2, we encounter arc $(7, 5)$ later than arc $(2, 3)$; so we select arc $(7, 5)$ as the leaving arc.

We show that the leaving arc rule guarantees that in the next spanning tree every node in the cycle W can send a positive amount of flow to the root node. Let

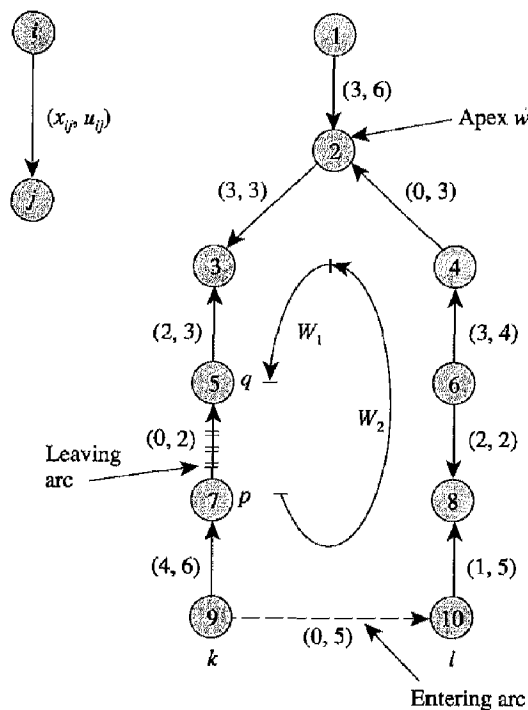


Figure 11.13 Selecting the leaving arc.

(p, q) be the arc selected by the leaving arc rule. Let W_1 be the segment of the cycle W between the apex w and arc (p, q) when we traverse the cycle along its orientation. Let $W_2 = W - W_1 - \{(p, q)\}$. Define the orientation of segments W_1 and W_2 as compatible with the orientation of W . See Figure 11.13 for an illustration of the segments W_1 and W_2 . We use the following property about the nodes in the segment W_2 .

Property 11.6. *Each node in the segment W_2 can send a positive amount of flow to the root in the next spanning tree.*

This observation follows from the fact that arc (p, q) is the last blocking arc in W ; consequently, no arc in W_2 is blocking and every node in this segment can send a positive amount of flow to the root via node w along the orientation of W_2 . Note that if the leaving arc does not satisfy the leaving arc rule, no node in the segment W_2 can send a positive amount of flow to the root; therefore, the next spanning tree will not be strongly feasible.

We next focus on the nodes contained in the segment W_1 .

Property 11.7. *Each node in the segment W_1 can send a positive amount of flow to the root in the next spanning tree.*

We prove this observation by considering two cases. If the previous pivot was a nondegenerate pivot, the pivot augmented a positive amount of flow δ along the arcs in W_1 ; consequently, after the augmentation, every node in the segment W_1 can send a positive amount of flow back to the root opposite to the orientation of W_1 via the apex node w (each node can send at least δ units to the apex and then at least some of this flow to the root since the previous spanning tree was strongly feasible). If the previous pivot was a degenerate pivot, W_1 must be contained in the segment of W between node w and node k because the property of strong feasibility implies that every node on the path from node l to node w can send a positive amount of flow to the root before the pivot, and thus no arc on this path can be a blocking arc in a degenerate pivot. Now observe that before the pivot, every node in W_1 could send a positive amount of flow to the root, and therefore since the pivot does not change flow values, every node in W_1 must be able to send a positive amount of flow to the root after the pivot as well. This conclusion completes the proof that in the next spanning tree every node in the cycle W can send a positive amount of flow to the root node.

We next show that in the next spanning tree, nodes not belonging to the cycle W can also send a positive amount of flow to the root. In the previous spanning tree (before the augmentation), every node j could send a positive amount of flow to the root and if the tree path from node j does not pass through the cycle W , the same path is available to carry a positive amount of flow in the next spanning tree. If the tree path from node j does pass through the cycle W , the segment of this tree path to the cycle W is available to carry a positive amount of flow in the next spanning tree and once a positive amount of flow reaches the cycle W , then, as shown earlier, we can send it (or some of it) to the root node. This conclusion completes the proof that the next spanning tree is strongly feasible.

We now establish the finiteness of the network simplex algorithm. Since we have previously shown that each nondegenerate pivot strictly decreases the objective function value, the number of nondegenerate pivots is finite. The algorithm can, however, also perform degenerate pivots. We will show that the number of successive degenerate pivots between any two nondegenerate pivots is finitely bounded. Suppose that arc (k, l) enters the spanning tree at its lower bound and in doing so it defines a degenerate pivot. In this case, the leaving arc belongs to the tree path from node k to the apex w . Now observe from Section 11.5 that node k lies in the subtree T_2 and the potentials of all nodes in T_2 change by an amount c_{kl}^T . Since $c_{kl}^T < 0$, this degenerate pivot strictly decreases the sum of all node potentials (which by our prior assumption is integral). Since no node potential can fall below $-nC$, the number of successive degenerate pivots is finite.

So far we have assumed that the entering arcs are always at their lower bounds. If the entering arc (k, l) is at its upper bound, we define the orientation of the cycle W as opposite to the orientation of arc (k, l) . The criteria for selecting the leaving arc remains unchanged—the leaving arc is the last blocking arc encountered in traversing W along its orientation starting at the apex w . In this case node l is contained in the subtree T_2 , and thus after the pivot, the potentials of all the nodes T_2 decrease by the amount $c_{kl}^T > 0$; consequently, the pivot again decreases the sum of the node potentials.

11.7 NETWORK SIMPLEX ALGORITHM FOR THE SHORTEST PATH PROBLEM

In this section we see how the network simplex algorithm specializes when applied to the shortest path problem. The resulting algorithm bears a close resemblance to the label-correcting algorithms discussed in Chapter 5. In this section we study the version of the shortest path problem in which we wish to determine shortest paths from a given source node s to all other nodes in a network. In other words, the problem is to send 1 unit of flow from the source to every other node along minimum cost paths. We can formulate this version of the shortest path problem as the following minimum cost flow model:

$$\text{Minimize } \sum_{(i,j) \in A} c_{ij}x_{ij}$$

subject to

$$\sum_{\{j:(i,j) \in A\}} x_{ij} - \sum_{\{j:(j,i) \in A\}} x_{ji} = \begin{cases} n-1 & \text{for } i = s \\ -1 & \text{for all } i \in N - \{s\} \end{cases}$$

$$x_{ij} \geq 0 \quad \text{for all } (i, j) \in A.$$

If the network contains a negative (cost)-directed cycle, this linear programming formulation would have an unbounded solution since we could send an infinite amount of flow along this cycle without violating any of the constraints (because the arc flows have no upper bounds). The network simplex algorithm we describe is capable of detecting the presence of a negative cycle, and if the network contains no such cycle, it determines the shortest path distances.

Like other minimum cost flow problems, the shortest path problem has a spanning tree solution. Because node s is the only source node and all the other nodes are demand nodes, the tree path from the source node to every other node is a directed path. This observation implies that the spanning tree must be a *directed out-tree rooted at node s* (see Figure 11.14 and the discussion in Section 4.3). As before, we store this tree using predecessor, depth, and thread indices. In a directed out-tree, every node other than the source has exactly one incoming arc but could have several outgoing arcs. Since each node except node s has unit demand, the flow of arc (i, j) is $|D(j)|$. [Recall that $D(j)$ is the set of descendants of node j in the spanning tree and, by definition, this set includes node j .] Therefore, every tree of the shortest path problem is nondegenerate, and consequently, the network simplex algorithm will never perform degenerate pivots.

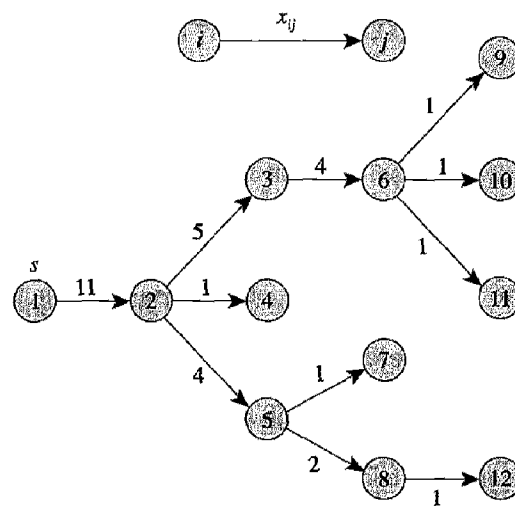


Figure 11.14 Directed out-tree rooted at the source.

Any spanning tree for the shortest path problem contains a unique directed path from node s to every other node. Let $P(k)$ denote the path from node s to node k . We obtain the node potentials corresponding to the tree T by setting $\pi(s) = 0$ and then using the equation $c_{ij} - \pi(i) + \pi(j) = 0$ for each arc $(i, j) \in T$ by fanning out from node s (see Figure 11.15). The directed out-tree property of the spanning tree implies that $\pi(k) = -\sum_{(i,j) \in P(k)} c_{ij}$. Thus $\pi(k)$ is the negative of the length of the path $P(k)$.

Since the variables in the minimum cost flow formulation of the shortest path problem have no upper bounds, every nontree arc is at its lower bound. The algorithm selects a nontree arc (k, l) with a negative reduced cost to introduce into the spanning tree. The addition of arc (k, l) to the tree creates a cycle which we orient in the same direction as arc (k, l) . Let w be the apex of this cycle. (See Figure 11.16 for an illustration.) In this cycle, every arc from node l to node w is a backward arc and every arc from node w to node k is a forward arc. Consequently, the leaving arc would lie in the segment from node l to node w . In fact, the leaving arc would be the arc $(\text{pred}(l), l)$ because this arc has the smallest flow value among all arcs in the segment from node l to node w . The algorithm would then increase the potentials of nodes in the subtree rooted at the node l by an amount $|c_{kl}^-|$, update the tree

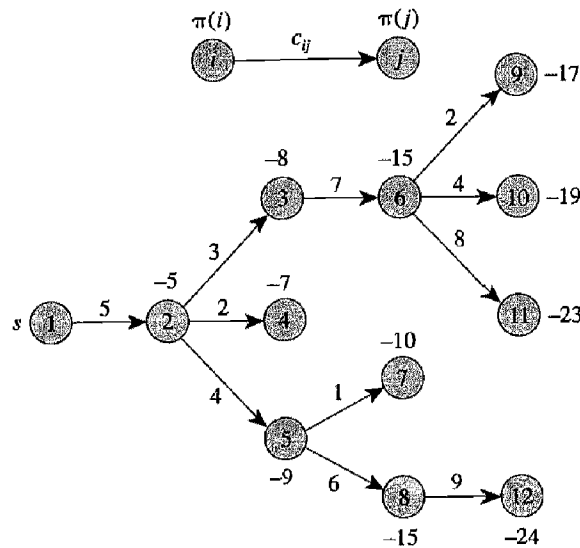


Figure 11.15 Computing node potentials.

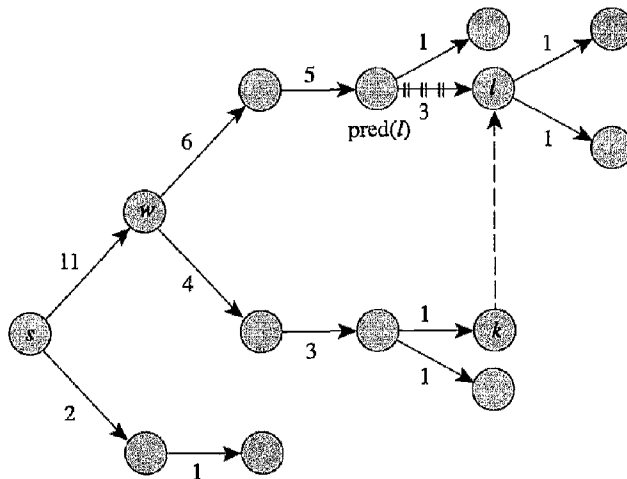


Figure 11.16 Selecting the leaving arc.

indices, and repeat the computations until all nontree arcs have nonnegative reduced costs. When the algorithm terminates, the final tree would be a shortest path tree (i.e., a tree in which the directed path from node s to every other node is a shortest path).

Recall that in implementing the network simplex algorithm for the minimum cost flow problem, we maintained flow values for all the arcs because we needed these values to identify the leaving arc. For the shortest path problem, however, we can determine the leaving arc without considering the flow values. If (k, l) is the entering arc, then $(\text{pred}(l), l)$ is the leaving arc. Thus the network simplex algorithm for the shortest path problem need not maintain arc flows. Moreover, updating of the tree indices is simpler for the shortest path problem.

The network simplex algorithm for the shortest path problem is similar to the label-correcting algorithms discussed in Section 5.3. Recall that a label-correcting algorithm maintains distance labels $d(i)$, searches for an arc satisfying the condition $d(j) > d(i) + c_{ij}$, and sets the distance label of node j equal to $d(i) + c_{ij}$. In the

network simplex algorithm, if we define $d(i) = -\pi(i)$, then $d(i)$ are the valid distance labels (i.e., they represent the length of some directed path from source to node i). At each iteration the network simplex algorithm selects an arc (i, j) with $c_{ij}^\pi < 0$. Observe that $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j) = c_{ij} + d(i) - d(j)$. Therefore, like a label-correcting algorithm, the network simplex algorithm selects an arc that satisfies the condition $d(j) > d(i) + c_{ij}$. The algorithm then increases the potential of every node in the subtree rooted at node j by an amount $|c_{ij}^\pi|$ which amounts to decreasing the distance label of all the nodes in the subtree rooted at node j by an amount $|c_{ij}^\pi|$. In this regard the network simplex algorithm differs from the label correcting algorithm: instead of updating one distance label at each step, it updates several of them.

If the network contains no negative cycle, the network simplex algorithm would terminate with a shortest path tree. When the network does contain a negative cycle, the algorithm would eventually encounter a situation like that depicted in Figure 11.17. This type of situation will occur only when the tail of the entering arc (k, l) belongs to $D(l)$, the set of descendants of node l . The network simplex algorithm can detect this situation easily without any significant increase in its computational effort: After introducing an arc (k, l) , the algorithm updates the potentials of all nodes in $D(l)$; at that time, it can check to see whether $k \in D(l)$, and if so, then terminate. In this case, tracing the predecessor indices would yield a negative cycle.

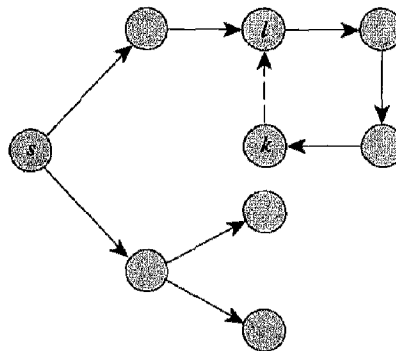


Figure 11.17 Detecting a negative cycle in the network.

The generic version of the network simplex algorithm for the shortest path problem runs in pseudopolynomial time. This result follows from the facts that (1) for each node i , $-nC \leq \pi(i) \leq nC$ (because the length of every directed path from s to node i lies between $-nC$ to nC), and (2) each iteration increases the value of at least one node potential. We can, however, develop special implementations that run in polynomial time. In the remainder of this section, in the exercises, and in the reference notes at the end of this chapter, we describe several polynomial-time implementations of the network simplex algorithm for the shortest path problem. These algorithms will solve the shortest path problem in $O(n^2m)$, $O(n^3)$, and $O(nm \log C)$ time. We obtain these polynomial-time algorithms by carefully selecting the entering arcs.

First eligible arc pivot rule. We have described this pivot rule in Section 11.5. The network simplex algorithm with this pivot rule bears a strong resemblance with the FIFO label-correcting algorithm that we described in Section 5.4. Recall

that the FIFO label-correcting algorithm examines the arc list in a wraparound fashion: If an arc (i, j) violates its optimality condition (i.e., it is eligible), the algorithm updates the distance label of node j . This order for examining the arcs ensures that after the k th pass, the algorithm has computed the shortest path distances to all those nodes that have a shortest path with k or fewer arcs. The network simplex algorithm with the first eligible arc pivot rule also examines the arc list in a wrap-around fashion, and if an arc (i, j) is eligible (i.e., violates its optimality condition), it updates the distances label of every node in $D(j)$, which also includes j . Consequently, this pivot rule will also, within k passes, determine shortest path distances to all those nodes that are connected to the source node s by a shortest path with k or fewer arcs. Consequently, the network simplex algorithm will perform at most n passes over the arc list. As a result, the algorithm will perform at most nm pivots and run in $O(n^2m)$ time. In Exercise 11.30 we discuss a modification of this algorithm that runs in $O(n^3)$ time.

Dantzig's pivot rule. This pivot rule selects the entering arc as an arc with the maximum violation. Let C denote the largest arc cost. We will show that the network simplex algorithm with this pivot rule performs $O(n^2 \log(nC))$ pivots and so runs in $O(n^2m \log(nC))$ time.

Scaled pivot rule. This pivot is a scaled variant of Dantzig's pivot rule. In this pivot rule we perform a number of scaling phases with varying values of a scaling parameter Δ . Initially, we let $\Delta = 2^{\lceil \log C \rceil}$ (i.e., we set Δ equal to the smallest power of 2 greater than or equal to C) and pivot in any nontree arc with a violation of at least $\Delta/2$. When no arc has a violation of at least $\Delta/2$, we replace Δ by $\Delta/2$ and repeat the steps. We terminate the algorithm when $\Delta < 1$.

We now show that the network simplex algorithm with the scaled pivot rule solves the shortest path problem in polynomial time. It is easy to verify that Dantzig's pivot rule is a special case of scaled pivot rule, so this result also shows that when implemented with Dantzig's pivot rule, the network simplex algorithm requires polynomial time.

We call the sequence of iterations for which Δ remains unchanged as the Δ -scaling phase. Let π denote the set of node potentials at the beginning of a Δ -scaling phase. Moreover, let $P^*(p)$ denote a shortest path from node s to node p and let $\pi^*(p) = -\sum_{(i,j) \in P^*} c_{ij}$ denote the optimal node potential of node p . Our analysis of the scaled pivot rule uses the following lemma:

Lemma 11.8. *If π denotes the current node potentials at the beginning of the Δ -scaling phase, then $\pi^*(p) - \pi(p) \leq 2n\Delta$ for each node p .*

Proof. In the first scaling phase, $\Delta \geq C$ and the lemma follows from the facts that $-nC$ and nC are the lower and upper bounds on any node potentials (why?). Consider next any subsequent scaling phase. Property 9.2 implies that

$$\sum_{(i,j) \in P^*(k)} c_{ij}^{\pi} = \sum_{(i,j) \in P^*(k)} c_{ij} - \pi(s) + \pi(p) = \pi(p) - \pi^*(p). \quad (11.4)$$

Since Δ is an upper bound on the maximum arc violation at the beginning of the Δ -scaling phase (except the first one), $c_{ij}^{\pi} \geq -\Delta$ for every arc $(i, j) \in A$. Sub-

stituting this inequality in (11.4), we obtain

$$\pi(p) - \pi^*(p) \geq -\Delta |P^*(p)| \geq -n\Delta,$$

which implies the conclusion of the lemma.

Now consider the potential function $\Phi = \sum_{p \in N} (\pi^*(p) - \pi(p))$. The preceding lemma shows that at the beginning of each scaling phase, Φ is at most $2n^2\Delta$. Now, recall from our previous discussion in this section that in each iteration, the network simplex algorithm increases at least one node potential by an amount equal to the violation of the entering arc. Since the entering arc has violation at least $\Delta/2$, at least one node potential increases by $\Delta/2$ units, causing Φ to decrease by at least $\Delta/2$ units. Since no node potential ever decreases, the algorithm can perform at most $4n^2$ iterations in this scaling phase. So, after at most $4n^2$ iterations, either the algorithm will obtain an optimal solution or will complete the scaling phase. Since the algorithm performs $O(\log C)$ scaling phases, it will perform $O(n^2 \log C)$ iterations and so require $O(n^2 m \log C)$ time. It is, however, possible to implement this algorithm in $O(nm \log C)$ time; the reference notes provide a reference for this result.

11.8 NETWORK SIMPLEX ALGORITHM FOR THE MAXIMUM FLOW PROBLEM

In this section we describe another specialization of the network simplex algorithm: its implementation for solving the maximum flow problem. The resulting algorithm is essentially an augmenting path algorithm, so it provides an alternative proof of the max-flow min-cut theorem we discussed in Section 6.5.

As we have noted before, we can view the maximum flow problem as a particular version of the minimum cost flow problem, obtained by introducing an additional arc (t, s) with cost coefficient $c_{ts} = -1$ and an upper bound $u_{ts} = \infty$, and by setting $c_{ij} = 0$ for all the original arcs (i, j) in A . To simplify our notation, we henceforth assume that A represents the set $A \cup \{(t, s)\}$. The resulting formulation is to

$$\text{Minimize} \quad -x_{ts}$$

subject to

$$\begin{aligned} \sum_{\{j: (i,j) \in A\}} x_{ij} - \sum_{\{j: (j,i) \in A\}} x_{ji} &= 0 \quad \text{for all } i \in N, \\ 0 \leq x_{ij} &\leq u_{ij} \quad \text{for all } (i, j) \in A. \end{aligned}$$

Observe that minimizing $-x_{ts}$ is equivalent to maximizing x_{ts} , which is equivalent to maximizing the net flow sent from the source to the sink, since this flow returns to the source via arc (t, s) . This observation explains why the inflow equals the outflow for every node in the network, including the source and the sink nodes.

Note that in any feasible spanning tree solution that carries a positive amount of flow from the source to the sink (i.e., $x_{ts} > 0$), arc (t, s) must be in the spanning tree. Consequently, the spanning tree for the maximum flow problem consists of

two subtrees of G joined by the arc (t, s) (see Figure 11.18). Let T_s and T_t denote the subtrees containing nodes s and t .

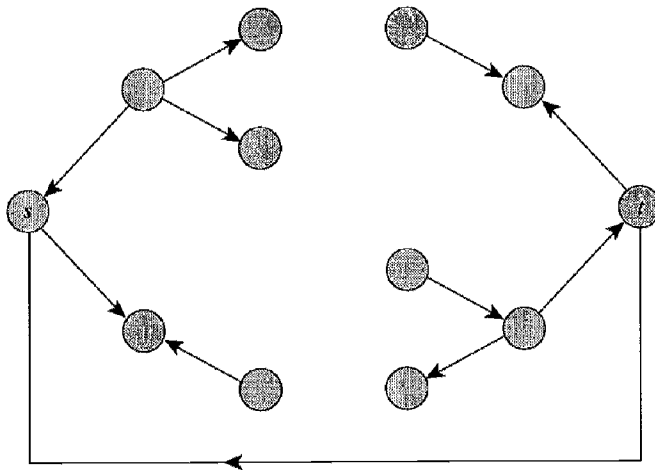


Figure 11.18 Spanning tree for the maximum flow problem.

We obtain node potentials corresponding to a feasible spanning tree of the maximum flow problem as follows. Since we can set one node potential arbitrarily, let $\pi(t) = 0$. Furthermore, since the reduced cost of arc (t, s) must be zero, $0 = c_{ts}^{\pi} = c_{ts} - \pi(t) + \pi(s) = -1 + \pi(s)$, which implies that $\pi(s) = 1$. Since (1) the reduced cost of every arc in T_s and T_t must be zero, and (2) the costs of these arcs are also zero, the node potentials have the following values: $\pi(i) = 1$ for every node $i \in T_s$, and $\pi(i) = 0$ for every node $i \in T_t$.

Notice that every spanning tree solution of the maximum flow problem defines an s - t cut $[S, \bar{S}]$ in the original network obtained by setting $S = T_s$ and $\bar{S} = T_t$. Each arc in this cut is a nontree arc; its flow has value zero or equals the arc's capacity. For every forward arc (i, j) in the cut, $c_{ij}^{\pi} = -1$, and for every backward arc (i, j) in the cut, $c_{ij}^{\pi} = 1$. Moreover, for every arc (i, j) not in the cut, $c_{ij}^{\pi} = 0$. Consequently, if every forward arc in the cut has a flow value equal to the arc's capacity and every backward arc has zero flow, this spanning tree solution satisfies the optimality conditions (11.1), and therefore it must be optimal. On the other hand, if in the current spanning tree solution, some forward arc in the cut has a flow of value zero or the flow on some backward arc equals the arc's capacity, all these arcs have a violation of 1 unit. Therefore, we can select any of these arcs to enter the spanning tree. Suppose that we select arc (k, l) . Introducing this arc into the tree creates a cycle that contains arc (t, s) as a forward arc (see Figure 11.19). The algorithm augments the maximum possible flow in this cycle and identifies a blocking arc. Dropping this arc again creates two subtrees joined by the arc (t, s) . This new tree constitutes a spanning tree for the next iteration.

Notice that this algorithm is an augmenting path algorithm: The tree structure permits us to determine the path from the source to the sink very easily. In this sense the network simplex algorithm has an advantage over other types of augmenting path algorithms for the maximum flow problem. As a compensating factor, however, due to degeneracy, the network simplex algorithm might not send a positive amount of flow from the source to the sink in every iteration.

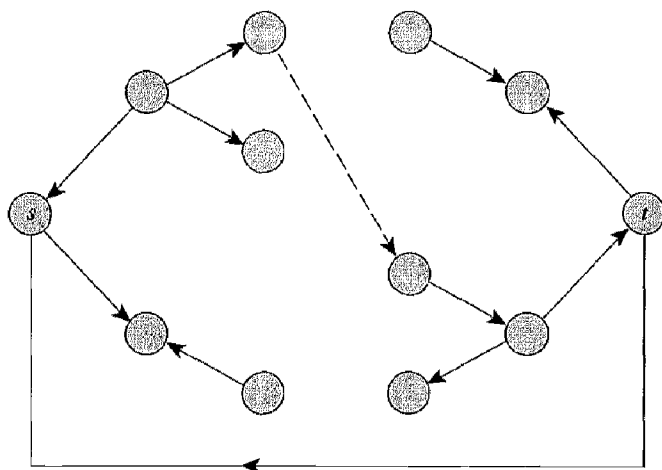


Figure 11.19 Forming a cycle.

The network simplex algorithm for the maximum flow problem gives another proof of the max-flow min-cut theorem. The algorithm terminates when every forward arc in the cut is capacitated and every backward arc has a flow of value zero. This termination condition implies that the current maximum flow value equals the capacity of the s - t cut defined by the subtrees T_s and T_t , and thus the value of a maximum flow from node s to node t equals the capacity of the minimum s - t cut.

Just as the mechanics of the network simplex algorithm becomes simpler in the context of the maximum flow problem, so does the concept of a strongly feasible spanning tree. If we designate the sink as the root node, the definition of a strongly feasible spanning tree implies that we can send a positive amount of flow from every node in T_t to the sink node t without violating any of the flow bounds. Therefore, every arc in T_t whose flow value is zero must point toward the sink node t and every arc in T_t whose flow value equals the arc's upper bound must point away from node t . Moreover, the leaving arc criterion reduces to selecting a blocking arc in the pivot cycle that is farthest from the sink node when we traverse the cycle in the direction of arc (t, s) starting from node t . Each degenerate pivot selects an entering arc that is incident to some node in T_t . The preceding observation implies that each blocking arc must be an arc in T_s . Consequently, each degenerate pivot increases the size of T_t , so the algorithm can perform at most n consecutive degenerate pivots. We might note that the minimum cost flow problem does not satisfy this property: For the more general problem, the number of consecutive degenerate pivots can be exponentially large.

The preceding discussion shows that when implemented to maintain a strongly feasible spanning tree, the network simplex algorithm performs $O(n^2U)$ iterations for the maximum flow problem. This result follows from the fact that the number of nondegenerate pivots is at most nU , an upper bound on the maximum flow value. This bound on the number of iterations is nonpolynomial, so is not satisfactory from a theoretical perspective. Developing a polynomial-time network simplex algorithm for the maximum flow problem remained an open problem for quite some time. However, researchers have recently suggested an entering arc rule that performs only $O(nm)$ iterations and can be implemented to run in $O(n^2m)$ time. This rule

selects entering arcs so that the algorithm augments flow along shortest paths from the source to the sink. We provide a reference for this result in the reference notes.

11.9 RELATED NETWORK SIMPLEX ALGORITHMS

In this section we study two additional algorithms for the minimum cost flow problem—the *parametric network simplex algorithm* and the *dual network simplex algorithm*—that are close relatives of the network simplex algorithm. In contrast to the network simplex algorithm, which maintains a feasible solution at each intermediate step, both of these algorithms maintain an infeasible solution that satisfies the optimality conditions; they iteratively attempt to transform this solution into a feasible solution. The solution maintained by the parametric network simplex algorithm satisfies all of the problem constraints except the mass balance constraints at two specially designated nodes, s and t . The solution maintained by the dual network simplex algorithm satisfies all of the mass balance constraints but might violate the lower and upper bound constraints on some arc flows. Like the network simplex algorithm, both algorithms maintain a spanning tree at every step and perform all computations using the spanning tree.

Parametric Network Simplex Algorithm

For the sake of simplicity, we assume that the network has one supply node (the source s) and one demand node (the sink t). We incur no loss of generality in imposing this assumption because we can always transform a network with several supply and demand nodes into one with a single supply and a single demand node. The parametric network simplex algorithm starts with a solution for which $b'(s) = -b'(t) = 0$, and gradually increases $b'(s)$ and $-b'(t)$ until $b'(s) = b(s)$ and $b'(t) = b(t)$. Let T be a shortest path tree rooted at the source node s in the underlying network. The parametric network simplex algorithm starts with zero flow and with (T, L, U) with $L = A - T$ and $U = \emptyset$ as the initial spanning tree structure. Since, by Assumption 9.5, all the arc costs are nonnegative, the zero flow is an optimal flow provided that $b(s) = b(t) = 0$. Moreover, since T is a shortest path tree, the shortest path distances $d(\cdot)$ to the nodes satisfy the condition $d(j) = d(i) + c_{ij}$ for each $(i, j) \in T$, and $d(j) \leq d(i) + c_{ij}$ for each $(i, j) \notin T$. By setting $\pi(j) = -d(j)$ for each node j , these shortest path optimality conditions become the optimality conditions (11.1) of the initial spanning tree structure (T, L, U) .

Thus the parametric network simplex algorithm starts with an optimal solution of a minimum cost flow problem that violates the mass balance constraints only at the source and sink nodes. In the subsequent steps, the algorithm maintains optimality of the solution and attempts to satisfy the violated constraints by sending flow from node s to node t along tree arcs. The algorithm stops when it has sent the desired amount ($b(s) = -b(t)$) of flow.

In each iteration the algorithm performs the following computations. Let (T, L, U) be the spanning tree structure at some iteration. The spanning tree T contains a unique path P from node s to node t . The algorithm first determines the maximum amount of flow δ that can be sent from s to t along P while honoring the flow bounds

on the arcs. Let \bar{P} and \underline{P} denote the sets of forward and backward arcs in P . Then

$$\delta = \min[\min\{u_{ij} - x_{ij} : (i, j) \in \bar{P}\}, \min\{x_{ij} : (i, j) \in \underline{P}\}].$$

The algorithm either sends δ units of flow along P , or a smaller amount if it would be sufficient to satisfy the mass balance constraints at nodes s and t . As in the network simplex algorithm, all the tree arcs have zero reduced costs; therefore, sending additional flow along the tree path from node s to node t preserves the optimality of the solution. If the solution becomes feasible after the augmentation, the algorithm terminates. If the solution is still infeasible, the augmentation creates at least one *blocking arc* (i.e., an arc that prohibits us from sending additional flow from node s to node t). We select one such blocking arc, say (p, q) , as the *leaving arc* and replace it by some nontree arc (k, l) , called the *entering arc*, so that the next spanning tree both (1) satisfies the optimality condition, and (2) permits additional flow to be sent from node s to node t . We accomplish this transition from one tree to another by performing a *dual pivot*. Recall from Section 11.5 that a (primal) pivot first identifies an entering arc and then the leaving arc. In contrast, a dual pivot first selects the leaving arc and then identifies the entering arc.

We perform a dual pivot in the following manner. We first drop the leaving arc from the spanning tree. Doing so gives us two subtrees T_s and T_t , with $s \in T_s$ and $t \in T_t$. Let S and \bar{S} be the subsets of nodes spanned by these two subtrees. Clearly, the cut $[S, \bar{S}]$ is an s - t cut and the entering arc (k, l) must belong to $[S, \bar{S}]$ if the next solution is to be a spanning tree solution. As earlier, we let (\bar{S}, S) denote the set of forward arcs and (S, \bar{S}) the set of backward arcs in the cut $[S, \bar{S}]$. Each arc in the cut $[S, \bar{S}]$ is at its lower bound or at its upper bound. We define the set Q of *eligible arcs* as the set

$$Q = ((S, \bar{S}) \cap L) \cup ((\bar{S}, S) \cap U),$$

that is, the set of forward arcs at their lower bound and the set of backward arcs at their upper bound. Note that if we add a noneligible arc to the subtrees T_s and T_t , we cannot increase the flow from node s to node t along the new tree path joining these nodes (since the arc lies on the path and would be a forward arc at its upper bound or a backward arc at its lower bound). If we introduce an eligible arc, the new path from node s to node t might be able to carry a positive amount of flow. Next, notice that if $Q = \emptyset$, we can send no additional flow from node s to node t . In fact, the cut $[S, \bar{S}]$ has zero residual capacity and the current flow from node s to node t equals the maximum flow. If $b(s)$ is larger than this flow value, the minimum cost flow problem is infeasible. We now focus on situations in which $Q \neq \emptyset$. Notice that we cannot select an arbitrary eligible arc as the entering arc, because the new spanning tree must also satisfy the optimality condition. For each eligible arc (i, j) , we define a number θ_{ij} in the following manner:

$$\theta_{ij} = \begin{cases} c_{ij}^{\pi} & \text{if } (i, j) \in L, \\ -c_{ij}^{\pi} & \text{if } (i, j) \in U. \end{cases}$$

Since the spanning tree structure (T, L, U) satisfies the optimality condition (11.1), $\theta_{ij} \geq 0$ for every eligible arc (i, j) . Suppose that we select some eligible arc (k, l) as the entering arc. It is easy to see that adding the arc (k, l) to $T_s \cup T_t$ decreases the potential of each node in \bar{S} by θ_{kl} units (throughout the computations, we maintain

that the node potential of the source node s has value zero). This change in node potentials decreases the reduced cost of each arc in (S, \bar{S}) by θ_{kl} units and increases the reduced cost of each arc in (\bar{S}, S) by θ_{kl} units. We have four cases to consider.

Case 1. $(i, j) \in (S, \bar{S}) \cap L$

The reduced cost of the arc (i, j) becomes $c_{ij}^\pi - \theta_{kl}$. The arc will satisfy the optimality condition (11.1b) if $\theta_{kl} \leq c_{ij}^\pi = \theta_{ij}$.

Case 2. $(i, j) \in (S, \bar{S}) \cap U$

The reduced cost of the arc (i, j) becomes $c_{ij}^\pi - \theta_{kl}$. The arc will satisfy the optimality condition (11.1c) regardless of the value of θ_{kl} because $c_{ij}^\pi \leq 0$.

Case 3. $(i, j) \in (\bar{S}, S) \cap L$

The reduced cost of the arc (i, j) becomes $c_{ij}^\pi + \theta_{kl}$. The arc will satisfy the optimality condition (11.1b) regardless of the value of θ_{kl} because $c_{ij}^\pi \geq 0$.

Case 4. $(i, j) \in (\bar{S}, S) \cap U$

The reduced cost of the arc (i, j) becomes $c_{ij}^\pi + \theta_{kl}$. The arc will satisfy the optimality condition (11.1c) provided that $\theta_{kl} \leq -c_{ij}^\pi = \theta_{ij}$.

The preceding discussion implies that the new spanning tree structure will satisfy the optimality conditions provided that

$$\theta_{kl} \leq \theta_{ij} \text{ for each } (i, j) \in ((S, \bar{S}) \cap L) \cup ((\bar{S}, S) \cap U) \equiv Q.$$

Consequently, we select the entering arc (k, l) to be an eligible arc for which $\theta_{kl} = \min\{\theta_{ij} : (i, j) \in Q\}$. Adding the arc (k, l) to the subtrees T_s and T_t gives us a new spanning tree structure and completes an iteration. We refer to this dual pivot as *degenerate* if $\theta_{kl} = 0$, and as *nondegenerate* otherwise. We repeat this process until we have sent the desired amount of flow from node s to node t .

It is easy to implement the parametric network simplex algorithm so that it runs in pseudopolynomial time. In this implementation, if an augmentation creates several blocking arcs, we select the one closest to the source as the leaving arc. Using inductive arguments, it is possible to show that in this implementation, the subtree T_s will permit us to augment a positive amount of flow from node s to every other node in T_s along the tree path. Moreover, in each iteration, when the algorithm sends no additional flow from node s to node t , it adds at least one new node to T_s . Consequently, after at most n iterations, the algorithm will send a positive amount of flow from node s to node t . Therefore, the parametric network simplex algorithm will perform $O(nb(s))$ iterations.

To illustrate the parametric network simplex algorithm, let us consider the same example we used to illustrate the network simplex algorithm. Figure 11.20(a) gives the minimum cost flow problem if we choose $s = 1$ and $t = 6$. Figure 11.20(b) shows the tree of shortest paths. All the nontree arcs are at their lower bounds. In the first iteration, the algorithm augments the maximum possible flow from node 1 to node 6 along the tree path 1–2–5–6. This path permits us to send a maximum of $\delta = \min\{u_{12}, u_{25}, u_{56}\} = \min\{8, 2, 6\} = 2$ units of flow. Augmenting 2 units along the path creates the unique blocking arc $(2, 5)$. We drop arc $(2, 5)$ from the tree, creating

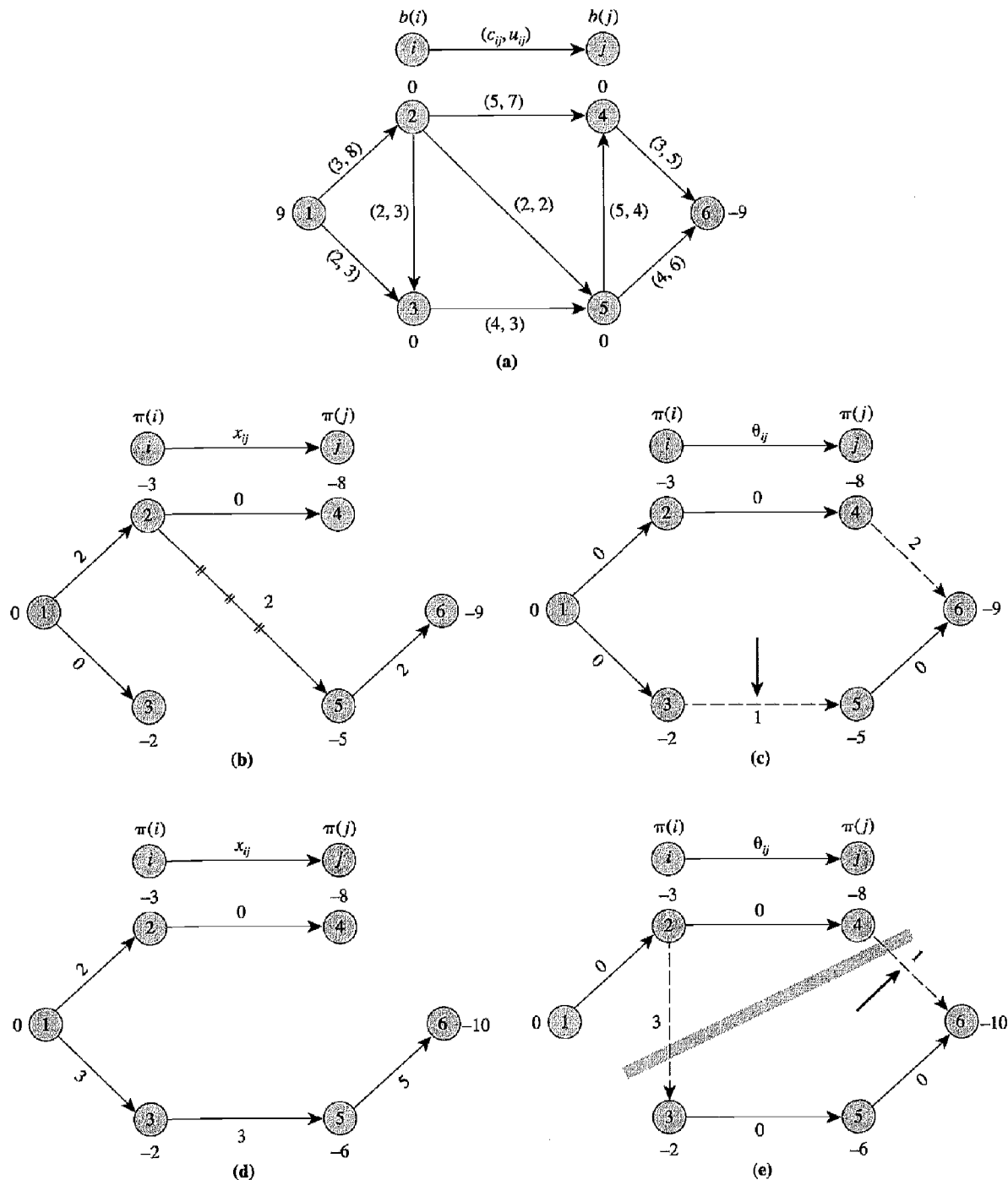


Figure 11.20 Illustrating the parametric network simplex algorithm.

the s - t cut $[S, \bar{S}]$ with $S = \{1, 2, 3, 4\}$ [see Figure 11.20(c)]. This cut contains two eligible arcs: arcs $(3, 5)$ and $(4, 6)$ with $\theta_{35} = 1$ and $\theta_{46} = 2$. We select arc $(3, 5)$ as the entering arc, creating the spanning tree shown in Figure 11.20(d). Notice that the potentials of the nodes 5 and 6 increase by 1 unit. In the new spanning tree, 1-

3–5–6 is the tree path from node 1 to node 6. We augment $\delta = \min\{u_{13}, u_{35}, u_{56} - x_{56}\} = \min\{3, 3, 6 - 2\} = 3$ units of flow along the path, creating two blocking arcs, (1, 3) and (3, 5). The arc (1, 3) is closer to the source and we select it as the leaving arc. As shown in Figure 11.20(e), the resulting s – t cut contains two eligible arcs, (2, 3) and (4, 6). Since $\theta_{46} < \theta_{23}$, we select (4, 6) as the entering arc. We leave the remaining steps of the algorithm as an exercise for the reader.

Notice the resemblance between the parametric network simplex algorithm and the successive shortest path algorithm that we discussed in Section 9.7. Both algorithms maintain the optimality conditions and gradually satisfy the mass balance constraints at the source and sink nodes. Both algorithms send flow along shortest paths from node s to node t . Whereas the successive shortest path algorithm does so by explicitly solving a shortest path problem, the parametric network simplex algorithm implicitly solves a shortest path problem. Indeed, the sequence of iterations that the parametric network simplex algorithm performs between two consecutive positive-flow iterations are essentially the steps of Dijkstra's algorithm for the shortest path problem.

Dual Network Simplex Algorithm

The dual network simplex algorithm maintains a solution that satisfies the mass balance constraints at all nodes, but that violates some of the lower and upper bounds imposed on the arc flows. The algorithm maintains a spanning tree structure (T, L, U) that satisfies the optimality conditions (11.1); the flow on the arcs in L and U are at their lower and upper bounds, but the flow on the tree arcs might not satisfy their flow bounds. We refer to a tree arc (i, j) as *feasible* if $0 \leq x_{ij} \leq u_{ij}$ and as *infeasible* otherwise. The algorithm attempts to make infeasible arcs feasible by sending flow along cycles; it terminates when the network contains no infeasible arc.

The dual network simplex algorithm performs a dual pivot at every iteration. Let (T, L, U) be the spanning tree structure at some iteration. In this solution some tree arcs might be infeasible. The algorithm selects any one of these arcs as the leaving arc. (Empirical evidence suggests that choosing an infeasible arc with the maximum violation of its flow bound generally results in a fewer number of iterations.) Suppose that we select the arc (p, q) as the leaving arc and $x_{pq} > u_{pq}$. We later address the case $x_{pq} < 0$. To make the flow on the arc (p, q) feasible, we must decrease the flow on this arc. We decrease the flow by introducing some nontree arc (k, l) that creates a unique cycle W containing arc (p, q) and augment enough flow along this cycle. Let us see which entering arc (k, l) would permit us to accomplish this objective.

If we drop the arc (p, q) from the spanning tree, we create two subtrees T_1 and T_2 , with $p \in T_1$ and $q \in T_2$. Let S and \bar{S} be the sets of nodes spanned by T_1 and T_2 . In addition, let (S, \bar{S}) and (\bar{S}, S) denote the sets of forward and backward arcs in the cut $[S, \bar{S}]$. Each arc in the cut $[S, \bar{S}]$, except the arc (p, q) , is at its lower or upper bound. Adding any arc (i, j) in $[S, \bar{S}]$ to T creates a unique cycle W that contains the arc (p, q) . Suppose that we define the orientation of the cycle W along the arc (i, j) if $(i, j) \in L$ and opposite to the arc (i, j) if $(i, j) \in U$. Each nontree arc in the cut $[S, \bar{S}]$ is (1) either a forward arc or a backward arc, and (2) either belongs to L or belongs to U . Consider any arc $(i, j) \in (S, \bar{S}) \cap L$. In this case, the orientation

of the cycle is along arc (i, j) ; consequently, arc (p, q) will be a backward arc in the cycle W and sending additional flow along the orientation of the cycle will decrease flow on the arc (p, q) [see Figure 11.21(a)]. Next, consider any arc $(i, j) \in (\bar{S}, S) \cap U$. In this case the orientation of the cycle is opposite to arc (i, j) ; therefore, sending additional flow along the orientation of the cycle again decreases flow on the arc (p, q) [see Figure 11.21(b)]. The reader can easily verify that in the other two cases when $(i, j) \in (S, \bar{S}) \cap U$ or $(i, j) \in (\bar{S}, S) \cap L$, increasing flow along the orientation of the cycle does not decrease flow on the arc (p, q) . Consequently, we define the set of *eligible arcs* as

$$Q = ((S, \bar{S}) \cap L) \cup ((\bar{S}, S) \cap U).$$

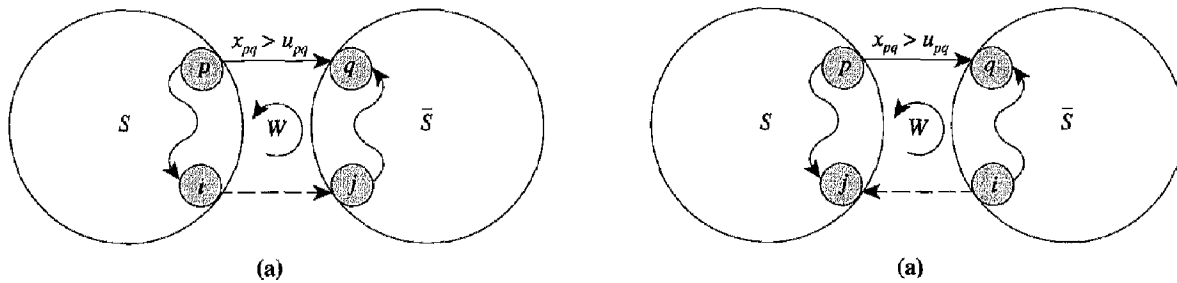


Figure 11.21 Identifying eligible arcs in the dual network simplex algorithm.

If $Q = \emptyset$, we cannot reduce the flow on arc (p, q) and the minimum cost flow problem is infeasible (see Exercise 11.37). If $Q \neq \emptyset$, we must select as the entering arc an eligible arc that would create a new spanning tree structure satisfying the optimality conditions. This step is similar to the same step in the parametric network simplex algorithm. We define a number θ_{ij} for each eligible arc (i, j) in the following manner:

$$\theta_{ij} = \begin{cases} c_{ij}^{\bar{}} & \text{if } (i, j) \in L, \\ -c_{ij}^{\bar{}} & \text{if } (i, j) \in U, \end{cases}$$

and select an arc (k, l) as the entering arc for which $\theta_{kl} = \min\{\theta_{ij} : (i, j) \in Q\}$. As before, we say this dual pivot is *degenerate* if $\theta_{kl} = 0$ and is *nondegenerate* otherwise. We augment $x_{pq} - u_{pq}$ units of flow along the cycle created by the arc (k, l) ; doing so decreases the flow on the arc (p, q) to value u_{pq} . Note that as a result of the augmentation, the arc (p, q) becomes feasible; other feasible arcs, however, might become infeasible. In the next spanning tree structure, the arc (k, l) replaces the arc (p, q) , and (p, q) becomes a nontree arc at its upper bound. Replacing the arc (p, q) by the arc (k, l) in the spanning tree decreases the potential of each node in \bar{S} by θ_{kl} units. (In the dual network simplex algorithm, the potential of node 1 might not always be zero.) As in our discussion of the parametric network simplex algorithm, it is possible to show that the new spanning tree structure satisfies the optimality conditions.

So far we have addressed situations in which the leaving arc (p, q) is infeasible because $x_{pq} > u_{pq}$. We now consider the case when $x_{pq} < 0$. In this instance, to make this arc feasible, we will increase its flow. The computations in this case are exactly the same as in the previous case except that we define the subtrees T_1 and

T_2 so that $p \in T_2$ and $q \in T_1$. We define the set of eligible arcs as $Q = ((S, \bar{S}) \cap L) \cup ((\bar{S}, S) \cap U)$ and select an eligible arc (k, l) with the minimum value of θ_{kl} as the entering arc. We augment $|x_{pq}|$ units of flow along the cycle created by the arc (k, l) ; doing so increases the flow on arc (p, q) to value zero. In the next spanning tree structure, arc (k, l) becomes a tree arc and (p, q) becomes a nontree arc at its lower bound.

Proving the finiteness of the dual network simplex algorithm is easy if each dual pivot is nondegenerate. As before, we assume that $x_{pq} > u_{pq}$ (a similar proof applies when $x_{pq} < 0$). In this case the entering arc (k, l) belongs to $(S, \bar{S}) \cap L$ or belongs to $(\bar{S}, S) \cap U$. In the former case, $c_{kl}^\pi > 0$ and the flow on the arc (k, l) increases by $(x_{pq} - u_{pq}) > 0$ units. In the latter case, $c_{kl}^\pi < 0$ and the flow on the arc decreases by $(x_{pq} - u_{pq}) > 0$ units. In either case, the cost of the flow increases by $c_{kl}^\pi(x_{pq} - u_{pq}) > 0$. Since mCU is an upper bound on the objective function value of the minimum cost flow problem and each nondegenerate pivot increases the cost by at least 1 unit, the dual network simplex algorithm will terminate finitely whenever every pivot is nondegenerate. In a degenerate pivot, the objective function value does not change because the entering arc (k, l) satisfies the condition $c_{kl}^\pi = 0$. In Exercise 11.38 we describe a dual perturbation technique that avoids the degenerate dual pivots altogether and yields a finite dual network simplex algorithm.

11.10 SENSITIVITY ANALYSIS

The purpose of sensitivity analysis is to determine changes in the optimal solution of the minimum cost flow problem resulting from changes in the data (supply/demand vector, capacity, or cost of any arc). In Section 9.11 we described methods for conducting sensitivity analysis using nonsimplex algorithms. In this section we describe network simplex based algorithms for performing sensitivity analysis.

Sensitivity analysis adopts the following basic approach. We first determine the effect of a given change in the data on the feasibility and optimality of the solution assuming that the spanning tree structure remains unchanged. If the change affects the optimality of the spanning tree structure, we perform (primal) pivots to achieve optimality. Whenever the change destroys the feasibility of the spanning tree structure, we perform dual pivots to achieve feasibility.

Let x^* denote an optimal solution of the minimum cost flow problem. Let (T^*, L^*, U^*) denote the corresponding spanning tree structure and π^* denote the corresponding node potentials. We first consider sensitivity analysis with respect to changes in the cost coefficients.

Cost Sensitivity Analysis

Suppose that the cost of an arc (p, q) increases by λ units. The analysis would be different when arc (p, q) is a tree or a nontree arc.

Case 1. Arc (p, q) is a nontree arc.

In this case, changing the cost of arc (p, q) does not change the node potentials of the current spanning tree structure. The modified reduced cost of arc (p, q) is $c_{pq}^{\pi^*} + \lambda$. If the modified reduced cost satisfies condition (11.1b) or (11.1c),

whichever is appropriate, the current spanning tree structure remains optimal. Otherwise, we reoptimize the solution using the network simplex algorithm with (T^*, L^*, U^*) as the starting spanning tree structure.

Case 2. Arc (p, q) is a tree arc.

In this case, changing the cost of arc (p, q) changes some node potentials. If arc (p, q) is an upward-pointing arc in the current spanning tree, potentials of all the nodes in $D(p)$ increase by λ , and if (p, q) is a downward-pointing arc, potentials of all the nodes in $D(q)$ decrease by λ . Note that these changes alter the reduced costs of those nontree arcs that belong to the cut $[D(q), \bar{D}(q)]$. If all nontree arcs still satisfy the optimality condition, the current spanning tree structure remains optimal; otherwise, we reoptimize the solution using the network simplex algorithm.

Supply/Demand Sensitivity Analysis

To study changes in the supply/demand vector, suppose that the supply/demand $b(k)$ of node k increases by λ and the supply/demand $b(l)$ of another node l decreases by λ . [Recall that since $\sum_{i \in N} b(i) = 0$, the supplies of two nodes must change simultaneously, by equal magnitudes and in opposite directions.] The mass balance constraints require that we must ship λ units of flow from node k to node l . Let P be the unique tree path from node k to node l . Let \bar{P} and \underline{P} , respectively, denote the sets of arcs in P that are along and opposite to the direction of the path. The maximum flow change δ_{ij} on an arc $(i, j) \in P$ that preserves the flow bounds is

$$\delta_{ij} = \begin{cases} u_{ij} - x_{ij} & \text{if } (i, j) \in \bar{P}, \\ x_{ij} & \text{if } (i, j) \in \underline{P}. \end{cases}$$

Let

$$\delta = \min\{\delta_{ij} : (i, j) \in P\}.$$

If $\lambda \leq \delta$, we send λ units of flow from node k to node l along the path P . The modified solution is feasible to the modified problem and since the modification in $b(i)$ does not affect the optimality of the solution, the resulting solution must be an optimal solution of the modified problem.

If $\lambda > \delta$, we cannot send λ units of flow from node k to node l along the arcs of the current spanning tree and preserve feasibility. In this case we send δ units of flow along P and reduce λ to $\lambda - \delta$. Let x' denote the updated flow. We next perform a dual pivot (as described in the preceding section) to obtain a new spanning tree that might allow additional flow to be sent from node k to node l along the tree path. In a dual pivot, we first decide on the leaving variable and then identify an entering variable. Let (p, q) be an arc in P that blocks us from sending additional flow from node k to node l . If $(p, q) \in \bar{P}$, then $x'_{pq} = u_{pq}$ and if $(p, q) \in \underline{P}$, then $x'_{pq} = 0$. We drop arc (p, q) from the spanning tree. Doing so partitions the set of nodes into two subtrees. Let S denote the subtree containing node k and \bar{S} denote the subtree containing node l . Now consider the cut $[S, \bar{S}]$. Since we wish to send additional flow through the cut $[S, \bar{S}]$, the arcs eligible to enter the tree would be the forward arcs in the cut at their lower bound or backward arcs at their upper bounds. If the

network contains no eligible arc, we can send no additional flow from node k to node l and the modified problem is infeasible. If the network does contain qualified arcs, then among these arcs, we select an arc, say (g, h) , whose reduced cost has the smallest magnitude. We introduce the arc (g, h) into the spanning tree and update the node potentials.

We then again try to send $\lambda' = \lambda - \delta$ units of flow from node k to node l on the tree path. If we succeed, we terminate; otherwise, we send the maximum possible flow and perform another dual pivot to obtain a new spanning tree structure. We repeat these computations until either we establish a feasible flow in the network or discover that the modified problem is infeasible.

Capacity Sensitivity Analysis

Finally, we consider sensitivity analysis with respect to arc capacities. Consider the analysis when the capacity of an arc (p, q) increases by λ units. (Exercise 11.40 considers the situation when an arc capacity decreases by λ units.) Whenever we increase the capacity of any arc, the previous optimal solution always remains feasible; to determine whether this solution remains optimal, we check the optimality conditions (11.1). If arc (p, q) is a tree arc or is a nontree arc at its lower bound, increasing u_{pq} by λ does not affect the optimality condition for that arc. If, however, arc (p, q) is a nontree arc at its upper bound and its capacity increases by λ units, the optimality condition (11.1c) dictates that we must increase the flow on the arc by λ units. Doing so creates an excess of λ units at node q and a deficit of λ units at node p . To achieve feasibility, we must send λ units from node q to node p . We accomplish this objective by using the method described earlier in our discussion of supply/demand sensitivity analysis.

11.11 RELATIONSHIP TO SIMPLEX METHOD

So far in this chapter, we have described the network simplex algorithm as a combinatorial algorithm and used combinatorial arguments to show that the algorithm correctly solves the minimum cost flow problem. This development has the advantage of highlighting the inherent combinatorial structure of the minimum cost flow problem and of the network simplex algorithm. The approach has the disadvantage, however, of not placing the network simplex method in the broader context of linear programming. To help to rectify this shortcoming, in this section we offer a linear programming interpretation of the network simplex algorithm. We show that the network simplex algorithm is indeed an adaptation of the well-known simplex method for general linear programs. Because the minimum cost flow problem is a highly structured linear programming problem, when we apply the simplex method to it, the resulting computations become considerably streamlined. In fact, we need not explicitly maintain the matrix representation (known as the simplex tableau) of the linear program and can perform all the computations directly on the network. As we will see, the resulting computations are exactly the same as those performed by the network simplex algorithm. Consequently, the network simplex algorithm is not a new minimum cost flow algorithm; instead, it is a special implementation of the

well-known simplex method that exploits the special structure of the minimum cost flow problem.

Our discussion in this section requires a basic understanding of the simplex method; Appendix C provides a brief review of this method. As we have noted before, the minimum cost flow problem is the following linear program:

Minimize cx

subject to

$$\mathcal{N}x = b,$$

$$0 \leq x \leq u.$$

The bounded variable simplex method for linear programming (or, simply, the simplex method) maintains a *basis structure* $(\mathbf{B}, \mathbf{L}, \mathbf{U})$ at every iteration and moves from one basis structure to another until it obtains an optimal basis structure. The set \mathbf{B} is the set of basic variables, and the sets \mathbf{L} and \mathbf{U} are the nonbasic variables at their lower and upper bounds. Following traditions in linear programming, we also refer to the variables in \mathbf{B} as a basis. Let \mathcal{B} , \mathcal{L} , and \mathcal{U} denote the sets of columns in \mathcal{N} corresponding to the variables in \mathbf{B} , \mathbf{L} , and \mathbf{U} . We refer to \mathcal{B} as a *basis matrix*. Our first result is a graph-theoretic characterization of the basis matrix.

Bases and Spanning Trees

We begin by establishing a one-to-one correspondence between bases of the minimum cost flow problem and spanning trees of G . One implication of this result is that the basis matrix is always lower triangular. The triangularity of the basis matrix is a key in achieving the efficiency of the network simplex algorithm.

We define the j th unit vector e_j as a column vector of size n consisting of all zeros except a 1 in the j th row. We let \mathcal{N}_{ij} denote the column of \mathcal{N} associated with the arc (i, j) . In Section 1.2 we show that $\mathcal{N}_{ij} = e_i - e_j$. The rows of \mathcal{N} are linearly dependent since summing all the rows yields the redundant constraint

$$0 = \sum_{i \in \mathcal{N}} b(i),$$

which is our assumption that the supplies/demands of all the nodes sum to zero. For convenience we henceforth assume that we have deleted the first row in \mathcal{N} (corresponding to node 1, which is treated as the root node). Thus \mathcal{N} has at most $n - 1$ independent rows. Since the number of linearly independent rows of a matrix is the same as the number of linearly independent columns, \mathcal{N} has at most $n - 1$ linearly independent columns. We show that the $n - 1$ columns associated with arcs of any spanning tree are linearly independent and thus define a basis matrix of the minimum cost flow problem.

Consider a spanning tree T . Let \mathcal{B} be the $(n - 1) \times (n - 1)$ matrix defined by the arcs in T . As an example, consider the spanning tree shown in Figure 11.22(a) which corresponds to the matrix \mathcal{B} shown in Figure 11.22(b). The first row in this matrix corresponds to the redundant row in \mathcal{N} and deleting this row yields an $(n - 1) \times (n - 1)$ square matrix. For the sake of clarity, however, we shall sometimes retain the first row. We order the rows and columns of \mathcal{B} in a certain specific

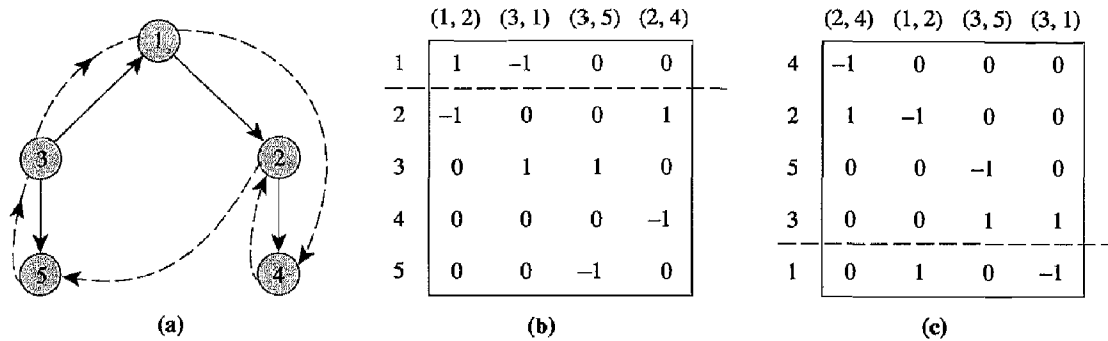


Figure 11.22 (a) Spanning tree and its reverse thread traversal; (b) basis matrix corresponding to the spanning tree; (c) basis matrix after rearranging the rows and columns.

manner. Doing so requires the reverse thread traversal of the nodes in the tree. Recall that a reverse thread traversal visits each node before visiting its predecessor. We order nodes and arcs in the following manner.

1. We order nodes of the tree in order of the reverse thread traversal. For our example, this order is 4–2–5–3–1 [see Figure 11.22(a)].
2. We order the tree arcs by visiting the nodes in order of the reverse thread traversal, and for each node i visited, we select the unique arc incident to it on the path to the root node. For our example, this order is (2, 4), (1, 2), (3, 5), and (3, 1).

We now arrange the rows and columns of \mathcal{B} as specified by the preceding node and arc orderings. Figure 11.22(c) shows the resulting matrix for our example. In this matrix, if we ignore the row corresponding to node 1, we have a lower triangular $(n - 1) \times (n - 1)$ matrix. The triangularity of the matrix is not specific to our example: The matrix would be triangular in general. It is easy to see why. Suppose that the reverse thread traversal selects node i at some step. Let $j = \text{pred}(i)$. Then either $(j, i) \in T$, or $(i, j) \in T$. Without any loss of generality, we assume that $(i, j) \in T$. The reverse thread traversal ensures that we have not visited node j so far. Consequently, the column corresponding to arc (i, j) will contain a +1 entry in the row r corresponding to node i , will contain all zero entries above this row, and will contain a -1 entry corresponding to node j below row r (because we will visit node j later). We have thus shown that this rearranged version of \mathcal{B} is a lower triangular matrix and that all of its diagonal elements are +1 or -1. We, therefore, have established the following result.

Theorem 11.9 (Triangularity Property). *The rows and columns of the node–arc incidence matrix of any spanning tree can be rearranged to be lower triangular.* ♦

The determinant of a lower triangular matrix is the product of its diagonal elements. Since each diagonal element in the matrix is ± 1 , the determinant is ± 1 . We now use the well-known fact from linear algebra that a set of $(n - 1)$ column

vectors, each of size $(n - 1)$, is linearly independent if and only if the matrix containing these vectors as columns has a nonzero determinant. This result shows that the columns corresponding to arcs of a spanning tree constitute a basis matrix of \mathcal{N} .

We now establish the converse result: Every basis matrix \mathcal{B} of \mathcal{N} defines a spanning tree. The fact that every basis matrix has the same number of columns implies that every basis matrix \mathcal{B} has $(n - 1)$ columns. These columns correspond to a subgraph G' of G having $(n - 1)$ arcs. Suppose that G' contains a cycle W . We assign any orientation to this cycle and consider the expression $\sum_{(i,j) \in W} (\pm 1) \mathcal{N}_{ij} = \sum_{(i,j) \in W} (\pm 1)(e_i - e_j)$; the leading coefficient of each term is $+1$ for those arcs aligned along the orientation of the cycle and is -1 for arcs aligned opposite to the orientation of the cycle. It is easy to verify that for each node j contained in the cycle, the unit vector e_j appears twice, once with a $+1$ sign and once with a -1 sign. Consequently, the preceding expression sums to zero, indicating that the columns corresponding to arcs of a cycle are linearly dependent. Since the columns of \mathcal{B} are linearly independent, G' must be an acyclic graph. Any acyclic graph on n nodes containing $(n - 1)$ arcs must be a spanning tree. So we have established the following theorem.

Theorem 11.10 (Basis Property). *Every spanning tree of G defines a basis of the minimum cost flow problem and, conversely, every basis of the minimum cost flow problem defines a spanning tree of G .* ♦

Implications of Triangularity

In the preceding discussion we showed that we can arrange every basis matrix of the minimum cost flow problem so that it is lower triangular and has an associated spanning tree. We now show that the triangularity of the basis matrix allows us to simplify the computations of the simplex method when applied to the minimum cost flow problem.

When applied to the minimum cost flow problem, the simplex method maintains a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$ at every step. Our preceding discussion implies that the arcs in the set \mathbf{B} constitute a spanning tree and the arcs in the set $\mathbf{L} \cup \mathbf{U}$ are nontree arcs. Therefore, this basis structure is no different from the spanning tree structure that the network simplex algorithm maintains. Moreover, the process of moving from one spanning tree structure to another corresponds to moving from one basis structure to another in the simplex method.

The simplex method performs the following operations:

1. Given a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$, determine the associated basic feasible solution.
2. Given a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$, determine the associated simplex multipliers π (or, dual variables).
3. Given a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$, check whether it is optimal, and if not, then determine an entering nonbasic variable x_{kl} .
4. Given a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$ and a nonbasic variable x_{kl} , determine the representation, $\bar{\mathcal{N}}_{kl}$, of the column \mathcal{N}_{kl} , corresponding to this variable in terms

of the basis matrix \mathcal{B} . We require this representation to perform the pivot operation while introducing the variable x_{kl} into the current basis.

We consider these simplex operations one by one.

Computing the Basic Feasible Solution

Given the basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$, the simplex method determines the associated basic feasible solution by solving the following system of equations:

$$\mathcal{B}x_{\mathbf{B}} = b - \mathcal{L}x_{\mathbf{L}} - \mathcal{U}u_{\mathbf{U}}. \quad (11.5)$$

In this expression, $x_{\mathbf{B}}$ denotes the set of basic variables, and $x_{\mathbf{L}}$ and $x_{\mathbf{U}}$ denote the sets of nonbasic variables at their lower and upper bounds. The simplex method sets each nonbasic variable in $x_{\mathbf{L}}$ to value zero, each nonbasic variable in $x_{\mathbf{U}}$ to its upper bound, and solves the resulting system of equations. Let $u_{\mathbf{U}}$ be the vector of upper bounds for variables in \mathbf{U} and let $b' = b - \mathcal{U}u_{\mathbf{U}}$. The simplex method solves the following system of equations:

$$\mathcal{B}x_{\mathbf{B}} = b'. \quad (11.6)$$

Let us see how can we solve (11.6) for the minimum cost flow problem. For simplicity of exposition, assume that $x_{\mathbf{B}} = (x_2, x_3, \dots, x_n)$. (Assume that the row corresponding to node 1 is the redundant row.) Since \mathcal{B} is a lower triangular matrix, the first row of \mathcal{B} has exactly one nonzero element corresponding to x_2 . Therefore, we can uniquely determine the value of x_2 . Since the coefficient of x_2 is ± 1 , the value of x_2 is integral. The second row of \mathcal{B} has at most two nonzero elements, corresponding to the variables x_2 and x_3 . Since we have already determined the value of x_2 , we can determine the value of x_3 uniquely. Continuing to solve successively for one variable at a time by this method of forward substitution, we can determine the entire vector $x_{\mathbf{B}}$. Since the nonzero coefficients in the basis matrix \mathcal{B} all have the value ± 1 , the only operations we perform are additions and subtractions, which preserve the integrality of the solution.

It is easy to see that the computations required to solve the system of equations $\mathcal{B}x_{\mathbf{B}} = b'$ are exactly same as those performed by the procedure *compute-flows* described in Section 11.4. Recall that the procedure first modifies the supply/demand vector b by setting the flows on the arcs in \mathbf{U} equal to their upper bounds. The modified supply/demand vector b' equals $b - \mathcal{U}u_{\mathbf{U}}$. Then the procedure examines the nodes in order of the reverse thread traversal and computes the flows on the arcs incident to these nodes. To put the matrix \mathcal{B} into a lower triangular form, we ordered its rows using the reverse thread traversal of the nodes. As a result, the procedure *compute-flows* computes flows on the arcs exactly in the same order as solving the system of equation $\mathcal{B}x_{\mathbf{B}} = b'$ by forward substitution.

Determining the Simplex Multipliers

The simplex algorithm determines the simplex multipliers π associated with a basis structure $(\mathbf{B}, \mathbf{L}, \mathbf{U})$ by solving the following system of equations:

$$\pi\mathcal{B} = c_{\mathbf{B}}. \quad (11.7)$$

In this expression, c_B is the vector consisting of cost coefficients of the variables in B . Assume, for simplicity of exposition, that $\pi = (\pi(2), \pi(3), \dots, \pi(n))$. Since \mathcal{B} is a lower triangular matrix, the last column of \mathcal{B} has exactly one nonzero element. Therefore, we can immediately determine $\pi(n)$. The second to last column of \mathcal{B} has at most two nonzero elements, corresponding to $\pi(n-1)$ and $\pi(n)$. Since we have already computed $\pi(n)$, we can easily compute $\pi(n-1)$, and so on. We can thus solve (11.7) by backward substitution and compute all the simplex multipliers by performing only additions and subtractions. Since we have arranged the rows of \mathcal{B} in the order of the reverse thread traversal of the nodes, and we determine simplex multipliers in the opposite order, we are, in fact, determining the simplex multipliers of nodes in the order dictated by the thread traversal. Recall from Section 11.4 that the procedure *compute-potentials* also examines nodes and computes the node potentials by visiting the nodes via the thread traversal. Consequently, the procedure *compute-potentials* is in fact solving the system of equations $\pi\mathcal{B} = c_B$ by backward substitution. Also, notice that the node potentials are the simplex multipliers maintained by the simplex method.

Optimality Testing

Given a basis structure (B, L, U) , the simplex method computes the simplex multipliers π , and then tests whether the basis structure satisfies the optimality conditions (11.1) (see Appendix C). As expressed in terms of the reduced costs c_{ij}^π , the optimality conditions are

$$c_{ij}^\pi = c_{ij} - \pi\mathcal{N}_{ij}, \quad \text{for each } (i, j) \in A.$$

For the minimum cost flow problem, $\mathcal{N}_{ij} = e_i - e_j$ and, therefore, $c_{ij}^\pi = c_{ij} - \pi(i) + \pi(j)$. Consequently, the reduced costs of the arcs as defined in the network simplex algorithm are the linear programming reduced costs and the optimality conditions (11.1) for the network simplex algorithm are the same as the linear programming optimality conditions (see Section C.5). The selection of the entering arc (k, l) in the network simplex algorithm corresponds to selecting the nonbasic variable x_{kl} as the entering variable. To simplify our subsequent exposition, we assume that the entering arc (k, l) is at its lower bound.

Representation of a Nonbasic Column

Once the simplex algorithm has identified a nonbasic variable x_{kl} to enter the basis, it next obtains the representation $\bar{\mathcal{N}}_{kl}$ of the column corresponding to x_{kl} with respect to the current basis matrix. We use this representation to determine the effect on the basic variables of assigning a value θ to x_{kl} , that is, to solve the system

$$x_B = \bar{b}' - \bar{\mathcal{N}}_{kl}\theta.$$

In this expression, $\bar{b}' = \mathcal{B}^{-1}b'$ and $\bar{\mathcal{N}}_{kl} = \mathcal{B}^{-1}\mathcal{N}_{kl}$. Observe that $-\bar{\mathcal{N}}_{kl}$ denotes the change in the values of basic variables as we increase the value of the entering nonbasic variable x_{kl} by 1 unit (i.e., set θ to value 1) and maintain all other nonbasic variables at their current lower and upper bounds. What is the graph-theoretic significance of $\bar{\mathcal{N}}_{kl}$?

The addition of arc (k, l) to the spanning tree T creates exactly one cycle, say W . Define the orientation of the cycle W to align with the orientation of the arc (k, l) . Let \overline{W} and \underline{W} denote the sets of forward and backward arcs in W . Observe that if we wish to increase the flow on arc (k, l) by 1 unit, keeping the flow on all other nontree arcs intact, then to satisfy the mass balance constraints we must augment 1 unit of flow along W . This change would increase the flow on arcs in \overline{W} by 1 unit and decrease the flow on arcs in \underline{W} by 1 unit. This discussion shows that the fundamental cycle W created by the nontree arc (k, l) defines the representation \overline{N}_{kl} in the following manner. All the basic variables corresponding to the arcs in \overline{W} have a coefficient of -1 in the column vector \overline{N}_{kl} , all the basic variables corresponding to the arcs in \underline{W} have a coefficient of $+1$, and all other basic variables have a coefficient of 0. This discussion also shows that in the network simplex algorithm, augmenting flow in the fundamental cycle created by the entering arc (k, l) and obtaining a new spanning tree solution corresponds to performing a pivot operation and obtaining a new basis structure in the simplex method.

To summarize, we have shown that the network simplex algorithm is the same as the simplex method applied to the minimum cost flow problem. The triangularity of the basis matrix permits us to apply the simplex method directly on the network without explicitly maintaining the simplex tableau. This possibility permits us to use the network structure to greatly improve the efficiency of the simplex method for solving the minimum cost flow problem.

In this section we have shown that the network simplex algorithm is an adaptation of the simplex method for solving general linear programs. A similar development would permit us to show that the parametric network simplex algorithm is an adaptation of the right-hand-side parametric algorithm of linear programming, and that the dual network simplex algorithm is an adaptation of the well-known dual simplex method for solving linear programs. We leave the details of these results as exercises (see Exercises 11.35 and 11.36).

11.12 UNIMODULARITY PROPERTY

In Section 11.4, using network flow algorithms, we established one of the fundamental results of network flows, the integrality property, stating that every minimum cost flow problem with integer supplies/demands and integer capacities has an integer optimal solution. The type of constructive proof that we used to establish this result has the obvious advantage of actually permitting us to compute integer optimal solutions. In that sense, constructive proofs have enormous value. However, constructive proofs do not always identify underlying structural (mathematical) reasons for explaining why results are true. These structural insights usually help in understanding a subject matter, and often suggest relationships between the subject matter and other problem domains or help to define potential limitations and generalization of the subject matter. In this section we briefly examine the structural properties of the integrality property, by providing an algebraic proof of this result. This discussion shows relationships between the integrality property and certain integrality results in linear programming.

Let \mathcal{A} be a $p \times q$ matrix with integer elements and p linearly independent rows (the matrix's rank is p). We say that the matrix \mathcal{A} is *unimodular* if the determinant

of every basis matrix \mathcal{B} of \mathcal{A} has value $+1$ or -1 [i.e., $\det(\mathcal{B}) = \pm 1$]. Recall from Appendix C that a $p \times p$ submatrix of \mathcal{A} is a basis matrix if its columns are linearly independent. The following classical result shows the relationship between unimodularity and the integer solvability of linear programs.

Theorem 11.11 (Unimodularity Theorem). *Let \mathcal{A} be an integer matrix with linearly independent rows. Then the following three conditions are equivalent:*

- (a) \mathcal{A} is unimodular.
- (b) Every basic feasible solution defined by the constraints $\mathcal{A}x = b$, $x \geq 0$, is integer for any integer vector b .
- (c) Every basis matrix \mathcal{B} of \mathcal{A} has an integer inverse \mathcal{B}^{-1} .

Proof. We prove the theorem by showing that (a) \Rightarrow (b), (b) \Rightarrow (c), and (c) \Rightarrow (a).

(a) \Rightarrow (b). Each basic feasible solution $x_{\mathcal{B}}$ has an associated basic matrix \mathcal{B} for which $\mathcal{B}x_{\mathcal{B}} = b$. By Cramer's rule, any component x_j of the solution $x_{\mathcal{B}}$ will be of the form

$$x_j = \frac{\det(\text{integer matrix})}{\det(\mathcal{B})}.$$

We obtain the integer matrix in this formula by replacing the j th column of \mathcal{B} with the vector b . Since, by assumption, \mathcal{A} is unimodular, $\det(\mathcal{B})$ is ± 1 , so x_j is integer.

(b) \Rightarrow (c). Let \mathcal{B} be a basis matrix of \mathcal{A} . Since \mathcal{B} has a nonzero determinant, its inverse \mathcal{B}^{-1} exists. Let e_j denote the j th unit vector (i.e., a vector with a 1 at the j th position and 0 elsewhere). Let $\mathcal{D} = \mathcal{B}^{-1}$ and \mathcal{D}_j denote the j th column of \mathcal{D} . We will show that the column vector \mathcal{D}_j is integer for each j whenever condition (b) holds. Select an integer vector α so that $\mathcal{D}_j + \alpha \geq 0$. Let $x = \mathcal{D}_j + \alpha$. Notice that

$$\mathcal{B}x = \mathcal{B}(\mathcal{D}_j + \alpha) = \mathcal{B}(\mathcal{B}^{-1}e_j + \alpha) = e_j + \mathcal{B}\alpha. \quad (11.8)$$

Multiplying the expression (11.8) by $\mathcal{D} = \mathcal{B}^{-1}$, we see that $x = \mathcal{D}_j + \alpha$. Since $e_j + \mathcal{B}\alpha$ is integer (by definition), condition (b) implies that $\mathcal{D}_j + \alpha$ is integer. Recalling that α is integer, we find that \mathcal{D}_j is also integer. This conclusion completes the proof of part (b).

(c) \Rightarrow (a). Let \mathcal{B} be a basis matrix of \mathcal{A} . By assumption, \mathcal{B} is an integer matrix, so $\det(\mathcal{B})$ is an integer. By condition (c), \mathcal{B}^{-1} is an integer matrix; consequently, $\det(\mathcal{B}^{-1})$ is also an integer. Since $\mathcal{B} \cdot \mathcal{B}^{-1} = I$ (i.e., an identity matrix), $\det(\mathcal{B}) \cdot \det(\mathcal{B}^{-1}) = 1$, which implies that $\det(\mathcal{B}) = \det(\mathcal{B}^{-1}) = \pm 1$. ♦

This result shows us when a linear program of the form minimize cx , subject to $\mathcal{A}x = b$, $x \geq 0$, has integer optimal solutions for *all* integer right-hand-side vectors b and for all cost vectors c . Network flow problems are the largest important class of models that satisfy this integrality property. To establish a formal connection between network flows and the results embodied in this theorem, we consider another noteworthy class of matrices.

Totally unimodular matrices are an important special subclass of unimodular

matrices. We say that a matrix \mathcal{A} is *totally unimodular* if each square submatrix of \mathcal{A} has determinant 0 or ± 1 . Every totally unimodular matrix \mathcal{A} is unimodular because each basis matrix \mathcal{B} must have determinant ± 1 (because the zero value of the determinant would imply the linear dependence of the columns of \mathcal{B}). However, a unimodular matrix need not be totally unimodular. Totally unimodular matrices are important, in large part, because the constraint matrices of the minimum cost flow problems are totally unimodular.

Theorem 11.12. *The node–arc incidence matrix \mathcal{N} of a directed network is totally unimodular.*

Proof. To prove the theorem, we need to show that every square submatrix \mathcal{F} of \mathcal{N} of size k has determinant 0, $+1$, or -1 . We establish this result by performing induction on k . Since each element of \mathcal{N} is 0, $+1$, or -1 , the theorem is true for $k = 1$. Now suppose that the theorem holds for some k . Let \mathcal{F} be any $(k + 1) \times (k + 1)$ submatrix of \mathcal{N} . The matrix \mathcal{F} satisfies exactly one of the three following possibilities: (1) \mathcal{F} contains a column with no nonzero element; (2) every column of \mathcal{F} has exactly two nonzero elements, in which case, one of these must be a $+1$ and the another a -1 ; and (3) some column \mathcal{F}_i has exactly one nonzero element, in, say, the i th row. In case (1) the determinant of \mathcal{F} is zero and the theorem holds. In case (2) summing all of the rows in \mathcal{F} yields the zero vector, implying that the rows in \mathcal{F} are linearly dependent and, consequently, $\det(\mathcal{F}) = 0$. In case (3) let \mathcal{F}' denote the submatrix of \mathcal{F} obtained by deleting the i th row and the i th column. Then $\det(\mathcal{F}) = \pm 1 \det(\mathcal{F}')$. By the induction hypothesis, $\det(\mathcal{F}')$ is 0, $+1$, or -1 , so $\det(\mathcal{F})$ is also 0, $+1$, or -1 . This conclusion establishes the theorem. ♦

This result, combined with Theorem 11.11, provides us with an algebraic proof of the integrality property of network flows: Network flow models have integer optimal solutions because every node–arc incidence matrix is totally unimodular and therefore unimodular. As we will see in later chapters, the constraint matrices for many extensions of the basic network flow problem, for example, generalized flows and multicommodity flows, are not unimodular. Therefore, we would not expect the optimal solutions of these models to be integer even when all of the underlying data are integer. Therefore, to find integer solutions to these problems, we need to rely on methods of integer programming. Although our development of the minimum cost flow problem has not stressed this point, one of the primary reasons that we are able to solve this problem so efficiently, and still obtain integer solutions, is because, as reflected by the integrality property, the basic feasible solutions of the linear programming formulation of this problem are integer whenever the underlying data are integer.

To close this section, we might note that the unimodularity properties provide us with a very strong result: any basic feasible solution is guaranteed to be integer-valued whenever the right-hand-side vector b is integer. It is possible, however, that basic feasible solutions to a linear program might be integer valued for a particular right-hand side even though they might be fractional for some other right-hand sides. We illustrate this possibility in Section 13.8 when we give an integer programming formulation of the minimal spanning tree problem.

11.13 SUMMARY

The network simplex algorithm is one of the most popular algorithms in practice for solving the minimum cost flow problem. This algorithm is an adaptation for the minimum cost flow problem of the well-known simplex method of linear programming. The linear programming basis of the minimum cost flow problem is a spanning tree. This property permits us to simplify the operations of the simplex method because we can perform all of its operations on the network itself, without maintaining the simplex tableau. Our development in this chapter does not require linear programming background because we have developed and proved the validity of the network simplex algorithm from first principles. Later in the chapter we showed the connection between the network simplex algorithm and the linear programming simplex method.

The development in this chapter relies on the fact that the minimum cost flow problem always has an optimal spanning tree solution. This result permits us to restrict our search for an optimal solution among spanning tree solutions. The network simplex algorithm maintains a spanning tree solution and successively transforms it into an improved spanning tree solution until it becomes optimal. At each iteration, the algorithm selects a nontree arc, introduces it into the current spanning tree, augments the maximum possible amount of flow in the resulting cycle, and drops a blocking arc from the spanning tree, yielding a new spanning tree solution. The algorithm is flexible in the sense that we can select the entering arc in a variety of ways and obtain algorithms with different worst-case and empirical attributes.

The network simplex algorithm does not necessarily terminate in a finite number of iterations unless we impose some additional restrictions on the choice of the entering and leaving arcs. We described a special type of spanning tree solution, called the *strongly feasible spanning tree solution*; when implemented in a way that maintains strongly feasible spanning tree solutions, the network simplex algorithm terminates finitely for any choice of the rule used for selecting the entering arc. We can maintain strongly feasible spanning tree solutions by selecting the leaving arc appropriately whenever several arcs qualify to be the leaving arc.

We also specialized the network simplex algorithm for the shortest path and maximum flow problems. When specialized for the shortest path problem, the algorithm maintains a directed out-tree rooted at the source node and iteratively modifies this tree until it becomes a tree of shortest paths. When we specialize the network simplex algorithm for the maximum flow problem, the algorithm maintains an s - t cut and selects an arc in this cut as the entering arc until the associated cut becomes a minimum cut.

The network simplex algorithm has two close relatives that might be quite useful in some circumstances: the parametric network simplex algorithm and the dual network simplex algorithm. The parametric network simplex algorithm maintains a spanning tree solution and parametrically increases the flow from a source node to a sink node until the algorithm has sent the desired amount of flow between these nodes. This algorithm is useful in situations in which we want to maximize the amount of flow to be sent from a source node to a sink node, subject to an upper bound on the cost of flow (see Exercise 10.25). The dual network simplex algorithm maintains a spanning tree solution in which spanning tree arcs do not necessarily satisfy the

flow bound constraints. The algorithm successively attempts to satisfy the flow bound constraints. The primary use of the dual network simplex algorithm has been for reoptimizing the minimum cost flow problem procedures for solving the minimum cost flow problem after we have changed the supply/demand or capacity data.

We also described methods for using the network simplex algorithm to conduct sensitivity analysis for the minimum cost flow problem with respect to the changes in costs, supplies/demands, and capacities. The resulting methods maintain a spanning tree solution and perform primal or dual pivots. Unlike the methods described in Section 9.11, these methods for conducting sensitivity analysis do not necessarily run in polynomial time (without further refinements). However, network simplex-based sensitivity analysis is excellent in practice.

The minimum cost flow problem always has an integer optimal solution; at the beginning of the chapter, we gave an algorithmic proof of this integrality property. We also examined the structural properties of the integrality property by providing an algebraic proof of this result. We showed that the constraint matrix of the minimum cost flow problem is totally unimodular and that, consequently, every basic feasible solution (or, equivalently, spanning tree solution) is an integer solution.

REFERENCE NOTES

Dantzig [1951] developed the network simplex algorithm for the uncapacitated transportation problem by specializing his linear programming simplex method. He proved the spanning tree property of the basis and the integrality property of the optimal solution. Later, his development of the upper bounding technique for linear programming led to an efficient specialization of the simplex method for the minimum cost flow problem. Dantzig's [1962] book discusses these topics.

The network simplex algorithm gained its current popularity in the early 1970s when the research community began to develop and test algorithms using efficient tree indices. Johnson [1966] suggested the first tree indices. Srinivasan and Thompson [1973], and Glover, Karney, Klingman, and Napier [1974] implemented these ideas; these investigations found the network simplex algorithm to be substantially faster than the existing codes that implemented the primal–dual and out-of-kilter algorithms. Subsequent research has focused on designing improved tree indices and determining the best pivot rule. The book by Kennington and Helgason [1980] describes a variety of tree indices and specifies procedures for updating them from iteration to iteration. The book by Bazaraa, Jarvis, and Sherali [1990] also describes a method for updating tree indices. The following papers describe a variety of pivot rules and the computational performance of the resulting algorithms: Glover, Karney, and Klingman [1974], Mulvey [1978], Bradley, Brown, and Graves [1977], Grigoriadis [1986], and Chang and Chen [1989]. The candidate list pivot rule that we describe in Section 11.5 is due to Mulvey [1978]. The reference notes of Chapter 9 contain information concerning the computational performance of the network simplex algorithm and other minimum cost flow algorithms.

Experience with solving large-scale minimum cost flow problems has shown that for certain classes of problems, more than 90% of the pivots in the network simplex algorithm can be degenerate. The strongly feasible spanning tree technique, proposed by Cunningham [1976] for the minimum cost flow problem, and indepen-

dently by Barr, Glover, and Klingman [1977] for the assignment problem, helps to reduce the number of degenerate steps in practice and ensures that the network simplex algorithm has a finite termination. Although the strongly feasible spanning tree technique prevents cycling during a sequence of consecutive degenerate pivots, the number of consecutive degenerate pivots can be exponential. This phenomenon is known as *stalling*. Cunningham [1979] and Goldfarb, Hao, and Kai [1990b] describe several antistalling pivot rules for the network simplex algorithm.

Researchers have attempted, with partial success, to develop polynomial-time implementations of the network simplex algorithm. Tarjan [1991] and Goldfarb and Hao [1988] have described polynomial-time implementations of a variant of the network simplex algorithm that permits pivots to increase value of the objective function. A monotone polynomial-time implementation, in which the value of the objective function is nonincreasing (as it does in any natural implementation), remains elusive to researchers.

Several FORTRAN codes of the network simplex algorithm are available in the public domain. These include (1) the RNET code developed by Grigoriadis and Hsu [1979], (2) the NETFLOW code developed by Kennington and Helgason [1980], and (3) a recent code by Chang and Chen [1989].

We next give selected references for several specific topics.

Shortest path problem. We have adapted the network simplex algorithm for the shortest path problem from Dantzig [1962]. Goldfarb, Hao, and Kai [1990a] and Ahuja and Orlin [1992a] developed the polynomial-time implementations of this algorithm that we have presented in Section 11.7. Additional polynomial-time implementations can be found in Orlin [1985] and Akgül [1985a].

Maximum flow problem. Fulkerson and Dantzig [1955] specialized the network simplex algorithm for the maximum flow problem. Goldfarb and Hao [1990] gave a polynomial-time implementation of this algorithm that performs at most nm pivots and runs in $O(n^2m)$ time; Goldberg, Grigoriadis, and Tarjan [1988] describe an $O(nm \log n)$ implementation of this algorithm.

Assignment problem. One popular implementation of the network simplex algorithm for the assignment problem is due to Barr, Glover, and Klingman [1977]. Roohy-Laleh [1980], Hung [1983], Orlin [1985], Akgül [1985b], and Ahuja and Orlin [1992a] have presented polynomial-time implementations of the network simplex algorithm for the assignment problem. Balinski [1986] and Goldfarb [1985] present polynomial-time dual network simplex algorithms for the assignment problem.

Parametric network simplex algorithm. Schmidt, Jensen, and Barnes [1982], and Ahuja, Batra, and Gupta [1984] are two sources for additional information on the parametric network simplex algorithm.

Dual network simplex algorithm. Ali, Padman, and Thiagarajan [1989] have described implementation details and computational results for the dual network simplex algorithm. Although no one has yet devised a (genuine) polynomial-time primal network simplex algorithm, Orlin [1984] and Plotkin and Tardos [1990]

have developed polynomial-time dual network simplex algorithms. The algorithm of Orlin [1984] is more efficient if capacities satisfy the similarity assumption; otherwise, the algorithm of Plotkin and Tardos [1990] is more efficient. The latter algorithm performs $O(m^2 \log n)$ pivots and runs in $O(m^3 \log n)$ time.

Sensitivity analysis. Srinivasan and Thompson [1972] have described parametric and sensitivity analysis for the transportation problem, which is similar to that for the minimum cost flow problem. Ali, Allen, Barr, and Kennington [1986] also discuss reoptimization procedures for the minimum cost flow problem.

Unimodularity. Hoffman and Kruskal [1956] first proved Theorem 11.11; the proof we have given is due to Veinott and Dantzig [1968]. The book by Schrijver [1986] presents an in-depth treatment of the unimodularity property and related topics.

EXERCISES

- 11.1. Nurse scheduling problem.** A hospital administrator needs to establish a staffing schedule for nurses that will meet the minimum daily requirements shown in Figure 11.23. Nurses reporting to the hospital wards for the first five shifts work for 8 consecutive hours, except nurses reporting for the last shift (2 A.M. to 6 A.M.), when they work for only 4 hours. The administrator wants to determine the minimal number of nurses to employ to ensure that a sufficient number of nurses are available for each period. Formulate this problem as a network flow problem.

Shift	1	2	3	4	5	6
Clock time	6 A.M. to 10 A.M.	10 A.M. to 2 P.M.	2 P.M. to 6 P.M.	6 P.M. to 10 P.M.	10 P.M. to 2 A.M.	2 A.M. to 6 A.M.
Minimum nurses required	70	80	50	60	40	30

Figure 11.23 Nurse scheduling problem.

- 11.2. Caterer problem.** As part of its food service, a caterer needs d_i napkins for each day of the upcoming week. He can buy new napkins at the price of α cents each or have his soiled napkins laundered. Two types of laundry service are available: regular and expedited. The regular laundry service requires two working days and costs β cents per napkin, and the expedited service requires one working day and costs γ cents per napkin ($\gamma > \beta$). The problem is to determine a purchasing and laundry policy that meets the demand at the minimum possible cost. Formulate this problem as a minimum costs flow problem. (*Hint:* Define a network on 15 nodes, 7 nodes corresponding to soiled napkins, 7 nodes corresponding to fresh napkins, and 1 node for the supply of fresh napkins.)
- 11.3. Project assignment.** In a new industry-funded academic program, each master's degree student is required to undertake a 6-month internship project at a company site. Since the projects are such an important component of the student's educational program

and vary considerably by company (e.g., by the problem and industry context) and by geography, each student would like to undertake a project of his or her liking. To assure that the project assignments are “fair,” the students and program administrators have decided to use an optimization approach: Each student ranks the available projects in order of increasing preference (lowest to highest). The objective is to assign students to projects to achieve the highest sum of total ranking of assigned projects. The project assignment has several constraints. Each student must work on exactly one project, and each project has an upper limit on the number of students it can accept. Each project must have a supervisor, drawn from a known pool of eligible faculty. Finally, each faculty member has bounds (upper and lower) on the number of projects that he or she can supervise. Formulate this problem as a minimum cost flow problem.

- 11.4. **Passenger routing.** United Airlines has six daily flights from Chicago to Washington. From 10 A.M. until 8 P.M., the flights depart every 2 hours. The first three flights have a capacity of 100 passengers and the last three flights can accommodate 150 passengers each. If overbooking results in insufficient room for a passenger on a scheduled flight, United can divert a passenger to a later flight. It compensates any passenger delayed by more than 2 hours from his or her regularly scheduled departure by paying \$200 plus \$20 for every hour of delay. United can always accommodate passengers delayed beyond the 8 P.M. flight on the 11 P.M. flight of another airline that always has a great deal of spare capacity. Suppose that at the start of a particular day the six United flights have 110, 160, 103, 149, 175, and 140 confirmed reservations. Show how to formulate the problem of determining the most economical passenger routing strategy as a minimum cost flow problem.
- 11.5. **Allocating receivers to transmitters** (Dantzig [1962]). An engine testing facility has four types of instruments: α_1 thermocouplers, α_2 pressure gauges, α_3 accelerometers, and α_4 thrust meters. Each instrument measures one type of engine characteristic and transmits its measurements over a separate communication channel. A set of receivers receive and record these data. The testing facility uses four types of receivers, each capable of recording one channel of information: β_1 cameras, β_2 oscilloscopes, β_3 instruments called “Idiots,” and β_4 instruments called “Hathaways.” The setup time of each receiver depends on the measurement instruments that are transmitting the data; let c_{ij} denote the setup time needed to prepare a receiver of type i to receive the information transmitted from any measurement taken by the j th instrument. The testing facility wants to find an allocation of receivers to transmitters that minimizes the total setup time. Formulate this problem as a network flow problem.
- 11.6. **Faculty–course assignment** (Mulvey [1979]). In 1973, the Graduate School of Management at UCLA revamped its M.B.A. curriculum. This change necessitated an increased centralization of the annual scheduling of faculty to courses. The large size of the problem (100 faculty, 500 courses, and three quarters) suggested that a mathematical model would be useful for determining an initial solution. The administration knows the courses to be taught in each of the three teaching quarters (fall, winter, and spring). Some courses can be taught in either of the two specified quarters; this information is available. A faculty member might not be available in all the quarters (due to leaves, sabbaticals, or other special circumstances) and when he is available he might be relieved from teaching some courses by using his project grants for “faculty offset time.” Suppose that the administration knows the quarters when a faculty member will be available and the total number of courses he will be teaching in those quarters. The school would like to maximize the preferences of the faculty for teaching the courses. The administration determines these preferences through an annual faculty questionnaire. The preference weights range from -2 to $+2$ and the administration occasionally revises the weights to reflect teaching ability and student inputs. Suggest a network model for determining a teaching schedule.
- 11.7. **Optimal rounding of a matrix** (Bacharach [1966], Cox and Ernst [1982]). In Application 6.3 we studied the problem of rounding the entries of a table to their nearest integers

while preserving the row and column sums of the matrix. We refer to any such rounding as a *consistent rounding*. Rounding off an element of the matrix introduces some error. If we round off an element a_{ij} to b_{ij} and $b_{ij} = \lfloor a_{ij} \rfloor$ or $b_{ij} = \lceil a_{ij} \rceil$, we measure the error as $(a_{ij} - b_{ij})^2$. Summing these terms for all the elements of the matrix gives us an error associated with any consistent rounding scheme. We say that a consistent rounding is an *optimal rounding* if the error associated with this rounding is as small as the error associated with any consistent rounding. Show how to determine an optimal rounding by solving a circulation problem. (*Hint*: Construct a network similar to the one used in Application 6.3. Define the arc costs appropriately.)

- 11.8. Describe an algorithm that either identifies p arc-disjoint directed paths from node s to node t or shows that the network does not contain any such set of paths. In the former case, show how to determine p arc-disjoint paths containing the fewest number of arcs. Suggest modifications of this algorithm to identify p node-disjoint directed paths from node s to node t containing the fewest number of arcs.
- 11.9. Show that a tree is a directed out-tree T rooted at node s if and only if every node in T except node s has indegree 1. State (but do not prove) an equivalent result for a directed in-tree.
- 11.10. Suppose that we permute the rows and columns of the node-arc incidence matrix N of a graph G . Is the modified matrix a node-arc incidence matrix of some graph G' ? If so, how are G' and G related?
- 11.11. Let T be a spanning tree of $G = (N, A)$. Every nontree arc (k, l) has an associated fundamental cycle which is the unique cycle in $T \cup \{(k, l)\}$. With respect to any arbitrary ordering of the arcs $(i_1, j_1), (i_2, j_2), \dots, (i_m, j_m)$, we define the *incidence vector* of any cycle W in G as an m -vector whose k th element is (1) 1, if (i_k, j_k) is a forward arc in W ; (2) -1 , if (i_k, j_k) is a backward arc in W ; and (3) 0, if $(i_k, j_k) \notin W$. Show how to express the incidence vector of any cycle W as a sum of incidence vectors of fundamental cycles.
- 11.12. Figure 11.24(b) gives a feasible solution of the minimum cost flow problem shown in Figure 11.24(a). Convert this solution into a spanning tree solution with the same or lower cost.

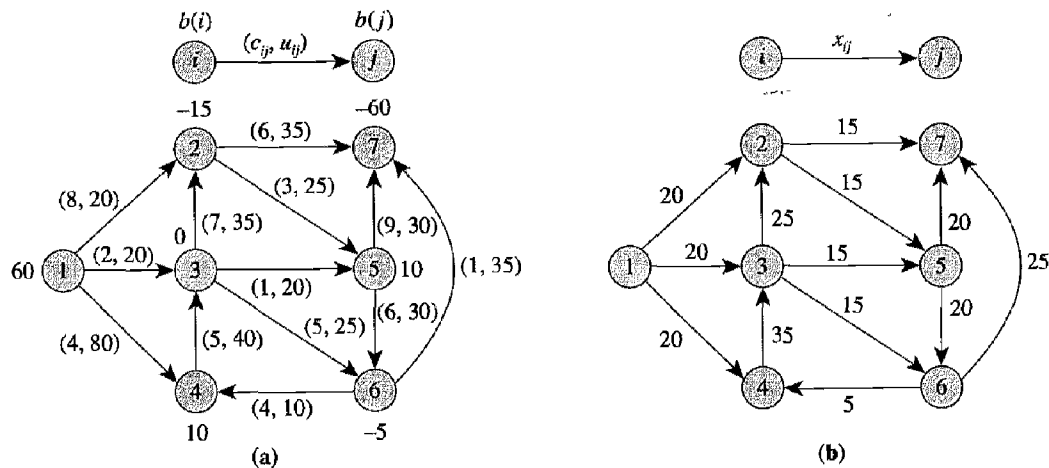


Figure 11.24 Example for Exercise 11.12: (a) problem data; (b) feasible solution.

- 11.13. Figure 11.25 specifies two spanning trees for the minimum cost flow problem shown in Figure 11.24(a). For Figure 11.25(a), compute the spanning tree solution assuming that all nontree arcs are at their lower bounds. For Figure 11.25(b), compute the spanning tree solution assuming that all nontree arcs are at their upper bounds.

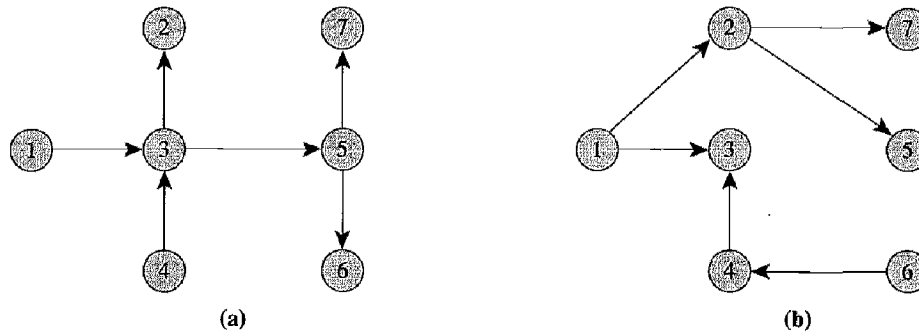


Figure 11.25 Two spanning trees of the network in Figure 11.24.

- 11.14. Assume that the spanning trees in Figure 11.25 have node 1 as their root. Specify the predecessor, depth, thread, and reverse thread indices of the nodes.
- 11.15. Compute the node potentials associated with the trees shown in Figure 11.25, which are the spanning trees of the minimum cost flow problem given in Figure 11.24(a). Verify that for each node j , the node potential $\pi(j)$ equals the length of the tree path from node j to the root.
- 11.16. Consider the minimum cost flow problem shown in Figure 11.26. Using the network simplex algorithm implemented with the first eligible pivot rule, find an optimal solution of this problem. Assume, as always, that arcs are arranged in the increasing order of their tail nodes, and for the same tail node, they are arranged in the increasing order of their head nodes. Use the following initial spanning tree structure: $T = \{(1, 2), (3, 2), (2, 5), (4, 5), (4, 6)\}$, $L = \{(3, 5)\}$, and $U = \{(1, 3), (2, 4), (5, 6)\}$.

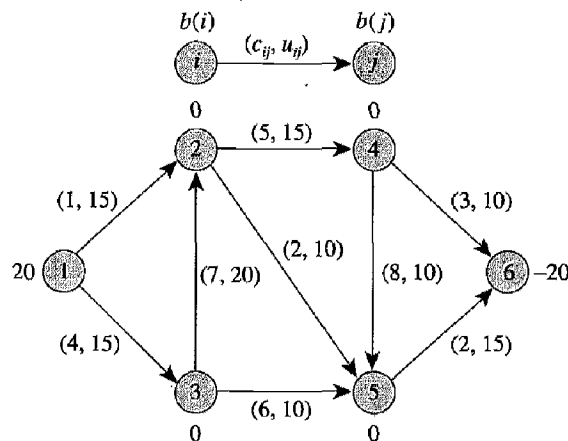


Figure 11.26 Example for Exercises 11.16 and 11.17.

- 11.17. Using the network simplex algorithm implemented with Dantzig's pivot rule, solve the minimum cost flow problem shown in Figure 11.26. Use the same initial spanning tree structure as used in Exercise 11.16.
- 11.18. In the procedure *compute-potentials*, we set $\pi(1) = 0$ and then compute other node potentials. Suppose, instead, that we set $\pi(1) = \alpha$ for some $\alpha > 0$ and then recompute all the node potentials. Show that all the node potentials increase by the amount α . Also show that this change does not affect the reduced cost of any arc.
- 11.19. Justify the procedure *compute-flows* for capacitated networks.
- 11.20. In the candidate list pivot rule, let *size* denote the maximum allowable size of the candidate list and *iter* denote the maximum number of minor iterations to be performed within a major iteration.

- (a) Specify values of size and iter so that the candidate list pivot rule reduces to Dantzig's pivot rule.
 - (b) Specify values of size and iter so that the candidate list pivot rule reduces to the first eligible arc pivot rule.
- 11.21. In Section 11.5 we showed how to find the apex of the pivot cycle W in $O(|W|)$ time using the predecessor and depth indices. Show that by using predecessor indices alone, you can find the apex of the pivot cycle in $O(|W|)$ time. (*Hint*: Do so by scanning at most $2|W|$ arcs.)
- 11.22. Given the predecessor indices of a spanning tree, describe an $O(n)$ time method for computing the thread and depth indices.
- 11.23. Describe methods for updating the predecessor and depth indices of the nodes when performing a pivot operation. Your method should require $O(n)$ time and should run faster than recomputing these indices from scratch.
- 11.24. Prove that in a spanning tree we can send a positive amount of flow from any node to the root without violating any flow bound if and only if every tree arc with zero flow is upward pointing and every tree arc at its upper bound is downward pointing.
- 11.25. Let $G(x)$ denote the residual network corresponding to a flow x . Show that a spanning tree T is a strongly feasible spanning tree if and only if for every node $i \in N - \{1\}$, $G(x)$ contains the arc $(i, \text{pred}(i))$.
- 11.26. **Primal perturbation.** In the minimum cost flow problem on a network G , suppose that we alter the supply/demand vector from value b to value $b + \epsilon$ for some vector ϵ . Let us refer to the modified problem as a *perturbed problem*. We consider the perturbation ϵ defined by $\epsilon(i) = 1/n$ for all $i = 2, 3, \dots, n$, and $\epsilon(1) = -(n - 1)/n$.
- (a) Let T be a spanning tree of G and let $D(j)$ denote the set of descendants of node j in T . Show that the perturbation decreases the flow on a downward-pointing arc (i, j) by the amount $|D(j)|/n$ and increases the flow on an upward-pointing arc (i, j) by the amount $|D(i)|/n$. Conclude that in a strongly feasible spanning tree solution, each arc flow is nonzero and is an integral multiple of $1/n$.
 - (b) Use the result in part (a) to show that the network simplex algorithm solves the perturbed problem in pseudopolynomial time irrespective of the pivot rule used for selecting entering arcs.
- 11.27. **Perturbation and strongly feasible solutions.** Let (T, L, U) be a feasible spanning tree structure of the minimum cost flow problem and let ϵ be a perturbation as defined in Exercise 11.26. Show that (T, L, U) is strongly feasible if and only if (T, L, U) remains feasible when we replace b by $b + \epsilon$. Use this equivalence to show that when implemented to maintain a strongly feasible basis, the network simplex algorithm runs in pseudopolynomial time irrespective of the pivot rule used for selecting entering arcs.
- 11.28. Apply the network simplex algorithm to the shortest path problem shown in Figure 11.27(a). Use a depth-first search tree with node 1 as the source node in the initial spanning tree solution and perform three iterations of the algorithm.
- 11.29. Apply the network simplex algorithm to the maximum flow problem shown in Figure 11.27(b). Use the following spanning tree as the initial spanning tree: a breadth-first search tree rooted at node 1 and spanning the nodes $N - \{t\}$ plus the arc (t, s) . Show three iterations of the algorithm.
- 11.30. Consider the application of the network simplex algorithm, implemented with the following pivot rule, for solving the shortest path problem. We examine all the nodes, one by one, in a wraparound fashion. Each time we examine a node i , we scan all incoming arcs at that node, and if the incoming arcs contain an eligible arc, we pivot in the arc with the maximum violation. We terminate when during an entire pass of the nodes, we find that no arc is eligible. Show when implemented with this pivot rule, the network simplex algorithm would perform $O(n^2)$ pivot operations and would run in $O(n^3)$ time. (*Hint*: The proof is similar to the proof of the first eligible arc pivot rule that we discussed in Section 11.7.)

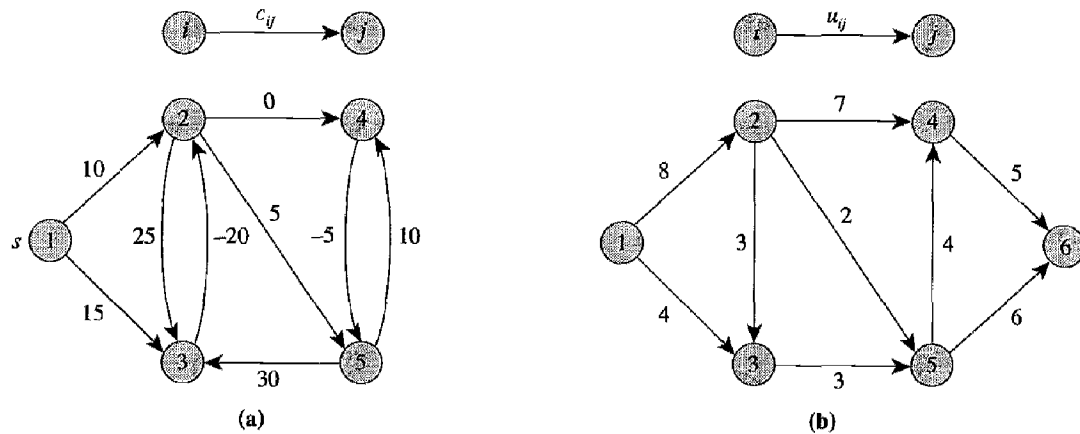


Figure 11.27 Examples for Exercises 11.28 and 11.29.

- 11.31.** The assignment problem, as formulated as a linear programming in (12.1), is a special case of the minimum cost flow problem. Show that every strongly feasible spanning tree of the assignment problem satisfies the following properties: (1) every downward-pointing arc carries unit flow; (2) every upward-pointing arc carries zero flow; and (3) every downward-pointing arc is the unique arc with flow equal to 1 emanating from node i .
- 11.32.** In a strongly feasible spanning tree of the assignment problem, a nontree arc (k, l) is a *downward arc* if node l is a descendant of node k . Show that when the network simplex algorithm, implemented to maintain strongly feasible spanning trees, is applied to the assignment problem, a pivot is nondegenerate if and only if the entering arc is a downward arc.
- 11.33.** Solve the minimum cost flow problem shown in Figure 11.26 by the parametric network simplex algorithm.
- 11.34.** Show how to solve the constrained maximum flow problem, as defined in Exercise 10.25, by a single application of the parametric network simplex algorithm.
- 11.35.** Show that the parametric network simplex algorithm described in Section 11.9 is an adaptation of the right-hand-side parametric simplex method of linear programming. (Consult any linear programming textbook for a review of the parametric simplex method of linear programming.)
- 11.36.** Show that the dual network simplex algorithm described in Section 11.9 is an adaptation of the dual simplex method of linear programming. (Consult any linear programming textbook for a review of the dual simplex method of linear programming.)
- 11.37.** At some point during its execution, the dual network simplex algorithm that we discussed in Section 11.9 might find that the set Q of eligible arcs is empty. In this case show that the minimum cost flow problem is infeasible. (*Hint:* Use the result in Exercise 6.43.)
- 11.38. Dual perturbation.** Suppose that we modify the cost vector c of a minimum cost flow problem on a network G in the following manner. After arranging the arcs in some order, we add $\frac{1}{2}$ to the cost of the first arc, $\frac{1}{4}$ to the cost of the second arc, $\frac{1}{8}$ to the cost of the third arc, and so on. We refer to the perturbed cost as c^* , and the minimum cost flow problem with the cost c^* as the *perturbed minimum cost flow problem*.
- (a) Show that if x^* is an optimal solution of the perturbed problem, x^* is also an optimal solution of the original problem. (*Hint:* Show that if $G(x^*)$ does not contain any negative cycle with cost c^* , it does not contain any negative cycle with cost c .)
- (b) Show that if we apply the dual network simplex algorithm to the perturbed problem, the reduced cost of each nontree arc is nonzero. Conclude that each dual

pivot in the algorithm will be nondegenerate and that the algorithm will terminate finitely. (*Hint*: Use the fact that the reduced cost of a nontree arc (k, l) is the cost of the fundamental cycle created by adding arc (k, l) to the spanning tree.)

- 11.39. In Exercise 9.24 we considered a numerical example concerning sensitivity analysis of a minimum cost flow problem. Solve the same problem using the simplex-based methods described in Section 11.10.
- 11.40. In Section 11.10 we described simplex-based procedures for reoptimizing a minimum cost flow solution when some cost coefficient c_{ij} increases or some flow bound u_{ij} decreases. Modify these procedures so that we can use them to handle situations in which (1) some c_{ij} decreases, or (2) some u_{ij} decreases.
- 11.41. Let \mathcal{B} denote the basis matrix associated with the columns of the spanning tree in Figure 11.25(a). Rearrange the rows and columns of \mathcal{B} so that it is lower triangular.
- 11.42. Let $G' = (N, A')$ be a subgraph of $G = (N, A)$ containing $|A'| = n - 1$ arcs. Let \mathcal{B}' be the square matrix defined by the columns of arcs in A' (where we delete one redundant row). Show that A' is a spanning tree of G if and only if the determinant of \mathcal{B}' is ± 1 .
- 11.43. **Computation of \mathcal{B}^{-1} .** In this exercise we discuss a combinatorial method for computing the inverse of a basis matrix \mathcal{B} of the minimum cost flow problem. (We assume that we have deleted a redundant row from \mathcal{B} .) By definition, $\mathcal{B}\mathcal{B}^{-1} = \mathcal{I}$, an identity matrix. Therefore, the j th column \mathcal{B}_j^{-1} of the inverse matrix \mathcal{B}^{-1} satisfies the condition $\mathcal{B}\mathcal{B}_j^{-1} = e_j$. Consequently, \mathcal{B}_j^{-1} is the unique solution x of the system of equations $\mathcal{B}x = e_j$. Assuming that we have deleted the row corresponding to node 1, x is the flow vector obtained from sending 1 unit of flow from node j to node 1 on the tree arcs corresponding to the basis. Use this result to compute \mathcal{B}^{-1} for the basis \mathcal{B} defined by the spanning trees shown in Figure 11.25(a).
- 11.44. Show that a matrix \mathcal{A} whose components are 0, +1, or -1 is totally unimodular if it satisfies both of the following conditions: (1) each column of \mathcal{A} contains at most two nonzero elements; and (2) the rows of \mathcal{A} can be partitioned into two subsets \mathcal{A}_1 and \mathcal{A}_2 so that the two nonzero entries in any column are in the same set of rows if they have different signs and are in different set of rows if they have the same sign.
- 11.45. Let \mathcal{N} be a totally unimodular matrix. Show that \mathcal{N}^T and $[\mathcal{N}, -\mathcal{N}]$ are also totally unimodular.
- 11.46. Show that a matrix \mathcal{N} is totally unimodular if and only if the matrix $[\mathcal{N}, \mathcal{I}]$ is unimodular.
- 11.47. Let T be a spanning tree of a directed network $G = (N, A)$ with node 1 as a designated root node. Let $d(i, j)$ denote the number of arcs on the tree path from node i to node j in T .
 - (a) For the given tree T , the *average depth* is $(\sum_{j \in N} d(1, j))/n$, and the *average cycle length* is $(\sum_{\text{nontree arcs } (i, j)} d(i, j) + 1)/(m - n + 1)$. Show that if G is a complete graph, the average cycle length is at most twice the average depth. Show that this relationship is not necessarily valid if the graph is not complete. (*Hint*: Use the fact that the length of the cycle created by adding the arc (i, j) to the tree is at most $d(1, i) + d(1, j) + 1$.)
 - (b) For a given tree T , let $D(j)$ denote the set of descendants of node j . The *average subtree size* of T is $(\sum_{j \in N} |D(j)|)/n$. Show that the average subtree size is 1 more than the average depth. (*Hint*: Let $E(j)$ denote the number of ancestors of node j in the tree T . First show that $\sum_{j \in N} |E(j)| = \sum_{j \in N} |D(j)|$.)
- 11.48. **Cost parametrization** (Srinivasan and Thompson [1972]). Suppose that we wish to solve a parametric minimum cost flow problem when the cost c_{ij} for each arc $(i, j) \in A$ is given by $c_{ij} = c_{ij}^0 + \lambda c_{ij}^*$ for some constants c_{ij}^0 and c_{ij}^* and we want to find an optimal solution for all values of the parameter λ in a given interval $[\alpha, \beta]$.
 - (a) Let (T, L, U) be an optimal spanning tree structure for the minimum cost flow problem for some value λ of the parameter. Let π^0 denote the node potentials for the tree T when c_{ij}^0 are the arc costs, and let π^* denote node potentials when

- c_{ij}^* are the arc costs in T (we can compute these potentials using the procedure *compute-potentials*). Show that $\pi^0 + \lambda\pi^*$ are the node potentials for the tree T when the arc costs are $c_{ij}^0 + \lambda c_{ij}^*$. Use this result to identify the largest value of λ , say $\bar{\lambda}$, for which (T, L, U) satisfies the optimality conditions.
- (b) Show that at $\lambda = \bar{\lambda}$, some nontree arc (k, l) satisfies its optimality condition as an equality and violates the optimality condition when $\lambda > \bar{\lambda}$. Show that if we perform the pivot operation with arc (k, l) as the entering arc, the new spanning tree structure also satisfies the optimality conditions at $\lambda = \bar{\lambda}$.
- (c) Use the results in parts (a) and (b) to solve the minimum cost flow problem for all values of the parameter λ in a given interval $[\alpha, \beta]$.
- 11.49. Supply/demand parametrization** (Srinivasan and Thompson [1972]). Suppose that we wish to solve a parametric minimum cost flow problem in which the supply/demand $b(i)$ of each node $i \in N$ is given by $b(i) = b^0(i) + \lambda b^*(i)$ for some constants $b^0(i)$ and $b^*(i)$ and we want to find an optimal solution for all values of the parameter λ in a given interval $[\alpha, \beta]$. We assume that $\sum_{i \in N} b^0(i) = \sum_{i \in N} b^*(i) = 0$.
- (a) Let (T, L, U) be an optimal spanning tree structure of the minimum cost flow problem for some value λ of the parameter. Let x_{ij}^0 and x_{ij}^* denote the flows on spanning tree arcs when b^0 and b^* are the supply/demand vectors (we can compute these flows using the procedure *compute-flows*). Show that $x_{ij}^0 + \lambda x_{ij}^*$ is the flow on the spanning tree arcs when $b^0 + \lambda b^*$ is the supply/demand vector. Use this result to identify the largest value of λ , say $\bar{\lambda}$, for which spanning tree arcs satisfy the flow bound constraints.
- (b) Show that at $\lambda = \bar{\lambda}$, some tree arc (p, q) satisfies one of its bounds (lower or upper bound) as an equality and violate its flow bound for $\lambda > \bar{\lambda}$. Show that if we perform a dual pivot (as described in Section 11.9) with arc (p, q) as the leaving arc, the new spanning tree structure also satisfies the optimality conditions at $\lambda = \bar{\lambda}$.
- (c) Use the results in parts (a) and (b) to solve the minimum cost flow problem for all values of the parameter λ in a given interval $[\alpha, \beta]$.
- 11.50. Capacity parametrization** (Srinivasan and Thompson [1972]). Consider a parametric minimum cost flow problem when the capacity u_{ij} of each arc $(i, j) \in A$ is given by $u_{ij} = u_{ij}^0 + \lambda u_{ij}^*$ for some constants u_{ij}^0 and u_{ij}^* . Describe an algorithm for solving the minimum cost flow problem for all values of the parameter λ in an interval $[\alpha, \beta]$. (*Hint: Let (T, L, U) be the basic structure at some state. Maintain the flow on each arc in the set U as the arc's upper flow bound (as a function of λ), determine the impact of this choice on the flows on the arcs in the spanning tree, and identify the maximum value of λ for which all the arc flows satisfy their flow bounds.*)
- 11.51. Constrained minimum cost flow problem.** The constrained minimum cost flow problem is a minimum cost flow problem with an additional constraint $\sum_{(i,j) \in A} d_{ij}x_{ij} \leq D$, called the *budget constraint*.
- (a) Show that the constrained minimum cost flow problem need not satisfy the integrality property (i.e., the problem need not have an integer optimal solution, even when all the data are integer).
- (b) For the constrained minimum cost flow problem, we say that a solution x is an *augmented tree solution* if some partition of the arc set A into the subsets $T \cup \{(p, q)\}$, L , and U satisfies the following two properties: (1) T is a spanning tree, and (2) by setting $x_{ij} = 0$ for each arc $(i, j) \in L$ and $x_{ij} = u_{ij}$ for each arc $(i, j) \in U$, we obtain a unique flow on the arcs in $T \cup \{(p, q)\}$ that satisfies the mass balance constraints and the budget constraint. Show that the constrained minimum cost flow problem always has an optimal augmented tree solution. Establish this result in two ways: (1) using a linear programming argument, and (2) using a combinatorial argument like the one we used in proving Theorem 11.2.